

This Page Is Inserted by IFW Operations  
and is not a part of the Official Record

## **BEST AVAILABLE IMAGES**

Defective images within this document are accurate representations of the original documents submitted by the applicant.

Defects in the images may include (but are not limited to):

- BLACK BORDERS
- TEXT CUT OFF AT TOP, BOTTOM OR SIDES
- FADED TEXT
- ILLEGIBLE TEXT
- SKEWED/SLANTED IMAGES
- COLORED PHOTOS
- BLACK OR VERY BLACK AND WHITE DARK PHOTOS
- GRAY SCALE DOCUMENTS

**IMAGES ARE BEST AVAILABLE COPY.**

**As rescanning documents *will not* correct images,  
please do not report the images to the  
Image Problem Mailbox.**

# PATENT ABSTRACTS OF JAPAN

(11)Publication number : 11-239158

(43)Date of publication of application : 31.08.1999

(51)Int.Cl.

H04L 12/28

H04L 1/22

H04Q 3/00

(21)Application number : 10-333325

(71)Applicant : THOMSON CSF

(22)Date of filing : 20.10.1998

(72)Inventor : DELATTRE MICHEL

BAVANT MARC

GUERIN DIDIER

HERAU PHILIPPE

(30)Priority

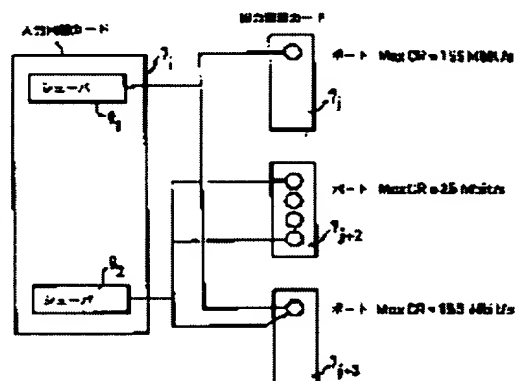
Priority number : 97 9713100 Priority date : 20.10.1997 Priority country : FR

## (54) FLOW CONTROL METHOD IN ATM SWITCH OF DISTRIBUTED CONSTITUTION

(57)Abstract:

**PROBLEM TO BE SOLVED:** To avoid the danger of causing congestion in a management module by distributing the (n) pieces of shapers adjusted as the function of a total average transmission speed dedicated to the flow of a VBRnrt category and the n-1 pieces of the shapers adjusted as the function of the available transmission speed of an output port dedicated to ABR and UBR category flow to respective input line cards.

**SOLUTION:** On the respective input line cards  $7i$ , the shapers  $91-9n$  are executed. One shaper is dedicated to the real time of the category VBRnrt (non real time variable transmission speed). The transmission speed of the shaper is adjusted as the function of the total average speed of the VBRnrt flow. The fixed number n-1 pieces of the shapers dedicated to the UBR (unspecified transmission speed)/ABR (available transmission speed) flow are adjusted as the function of the AvCR (available speed) of the output port. An operation range for featuring the value of the transmission speed achievable by the shaper is allocated to the respective UBR/ABR shapers identified by the number.



## LEGAL STATUS

[Date of request for examination]

[Date of sending the examiner's decision of

rejection]

[Kind of final disposal of application other than  
the examiner's decision of rejection or  
application converted registration]

[Date of final disposal for application]

[Patent number]

[Date of registration]

[Number of appeal against examiner's decision of  
rejection]

[Date of requesting appeal against examiner's  
decision of rejection]

[Date of extinction of right]

Copyright (C); 1998,2003 Japan Patent Office

(19) 日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11) 特許出願公開番号

特開平11-239158

(43) 公開日 平成11年(1999) 8月31日

(51) Int.Cl.<sup>6</sup>

識別記号

F I

H 0 4 L 12/28

H 0 4 L 11/20

G

1/22

1/22

H 0 4 Q 3/00

H 0 4 Q 3/00

審査請求 未請求 請求項の数 8 O L 外国語出願 (全 47 頁)

(21) 出願番号

特願平10-333325

(22) 出願日

平成10年(1998)10月20日

(31) 優先権主張番号

9 7 1 3 1 0 0

(32) 優先日

1997年10月20日

(33) 優先権主張国

フランス (F R)

(71) 出願人 591000827

トムソン・セーエスエフ

THOMSON-CSF

フランス国、75008・パリ、ブルパール・

オースマン・173

(72) 発明者 ミシェル デラトル

フランス国、92100 プーローニュ、

リュ ドゥ ベルヴュー、95番地

(72) 発明者 マルク バヴァン

フランス国、75019 パリ、リュ ク

ラヴェル、21番地

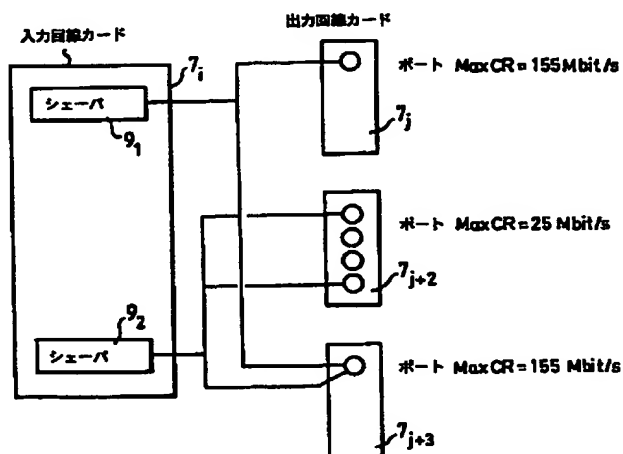
(74) 代理人 弁理士 山本 恵一

最終頁に続く

(54) 【発明の名称】 分散構成のATMスイッチでのフロー制御方法

(57) 【要約】

本方法は、各入力回線カード ( $7_i$ ) へ VBR  $n r t$  (非実時間可変伝送速度) カテゴリ・フロー専用の合計平均伝送速度の関数として調整された指定数  $n$  個のシェーパ ( $9_1 \dots 9_n$ ) と、出力ポートの利用可能伝送速度 (AVCR) の関数として調整された  $n-1$  個の別のシェーパを分散することからなる。ATM伝送ネットワークに適用される。



## 【特許請求の範囲】

【請求項1】 分散アーキテクチャを備え入力部に記憶を備えたATMスイッチ内の、各入力回線カードへVBR<sub>nrt</sub>（非実時間可変伝送速度）カテゴリ・フロー専用の合計平均伝送速度の関数として調整された指定数 $n$ 個のシェーパと、ABRおよびUBRカテゴリ・フロー専用の $n-1$ 個の別のシェーパを分散することからなるフロー制御の方法であって、それぞれのABRまたはUBR接続に根がその入力回線カードで構成され葉が出力ポートで構成される木を対応させ、この木をその葉の公称伝送速度の関数として選択されたシェーパに割り当て、シェーパがこのシェーパに割り当てられたすべての木の葉の利用可能伝送速度の関数として調整されることを特徴とする方法。

【請求項2】 それぞれのシェーパにPDUアプリケーション・フレーム・モードで動作する第1の待ち行列 $P_1$ とセル・モードで動作する第2の待ち行列 $P_2$ を割り当てる請求項1に記載の方法。

【請求項3】 第1の待ち行列 $P_1$ を多地点一地点間および多地点一地点間接続に予約する請求項2に記載の方法。

【請求項4】 それぞれのシェーパを動作範囲によって定義されたクラスに割り当て、それぞれのシェーパに関連付けられたそれぞれの木が少なくとも1つの共通の葉を備えている場合、この範囲のシェーパが別のシェーパと同じクラスに属する請求項1から3のいずれか一項に記載の方法。

【請求項5】 1つの同じクラスのシェーパに《オーダー・オブ・スピーキング（order of speaking）》を割り当て、ライト・トゥ・スピーク（right to speak）を、 $T=T_1+T_2+T_3$ になるように3つの基本時間間隔 $T_1$ 、 $T_2$ 、 $T_3$ に分割される、指定された時間間隔 $T$ に制限し、送信の先頭が間隔 $T_1$ から開始し送信の最後が間隔 $T_1+T_2$ に実行され、次の間隔 $T$ のターン・トゥ・スピーク（turn to speak）が間隔 $T_3$ でそれぞれのシェーパが表明した希望の関数としてPDUアプリケーション・フレームの送信を中断したシェーパのために決定される請求項4に記載の方法。

【請求項6】 周期性 $T$ を備えた《SYNC SIGNAL トークン》を一般同報通信でランク $i$ のシェーパが少なくともそれに関連付けられた1つの木を所有するすべての回線カードへ送信するための《マスタ》回線カードを指定し、ランク $i$ のシェーパが考慮されるクラスのターン・トゥ・スピーク（turn to speak）で $T_1+T_2$ 間隔の終了前に次の回線カードへの送信を完了した回線カードの1つによって《speech》トークンを送信し、時間間隔 $T_3$ の周期 $T$ の最後にSYNC SIGNAL トークンを受信した回線カードによって、後続の時間間隔 $T$ に送信を続行する必要を示す標識

を内部に記録するマスタ回線カードに向けて《collision》トークンを送信する請求項5に記載の方法。

【請求項7】  $VLA_i \rightarrow (VLA'_i, LG)$  から  $VLA_i \rightarrow (VLM, L1, L2, L3, \dots, Lq)$  への二地点間接続（ $A_i$ ）、ただし $VLA_i$ は入力幹線の $A_i$ に関連付けられた論理チャネル識別子、 $VLA'_i$ は管理モジュール内のこの同じ接続に関連付けられた論理チャネル識別子、 $LG$ は管理モジュールのスイッチ・ファブリック・ジャンクションの識別、 $VLM$ は一地点多地点間接続 $M$ の同報通信インデックス、 $L1$ 、 $L2$ 、 $\dots$ 、 $Lq$ は接続 $M$ が関連するすべてのスイッチ・ファブリック機能の識別を指定する、二地点間接続（ $A_i$ ）の入力翻訳を変更することでブロードキャスト・サーバのフレーム中継機能を回線カードのATMレイヤ内に移行して多地点間モードでの通信を実行する請求項6に記載の方法。

【請求項8】 そこから宛先のIPアドレスを抽出するためにそれぞれのPDUフレームの最初のセルを検査し、キャッシュ・テーブルを検索して論理チャネルと送信方向で構成されるペアを関連するIPアドレスの前に見つけ、PDUフレーム内のすべてのセルで得られる翻訳を用い、管理モジュール内に常駐するルーティング・エミュレーション機能から受信したルーティング情報によってキャッシュ・テーブルを更新して、ルーティング・エミュレーション装置のフレーム中継機能を回線カードのATMレイヤ内に移行して多地点間モードでの通信を実行し、所望のIPアドレスが見つからない場合にキャッシュの更新要求を管理モジュールへ送信する請求項6に記載の方法。

## 【発明の詳細な説明】

## 【0001】

【発明の属する技術分野】 本発明は、分散アーキテクチャおよび入力時の記憶を備えたATMスイッチ内の二地点間、一地点多地点間、多地点一地点間および多地点間の非リアルタイム接続のフロー制御の方法に関する。

## 【0002】

【従来の技術】 ATMすなわち非同期転送モード・ネットワークとして知られる通信ネットワークを用いて5バイトのヘッダと48バイトの本文からなるATMセルとして知られる固定長パケットの循環が可能になる。より詳細に言えば、ヘッダは送信元ユーザと宛先ユーザの間の経路上で検出するスイッチ内のセルのルーティングを可能にするVPI/VCI（すなわち仮想経路識別および仮想チャネル識別子）フィールドとして知られる論理チャネル識別子を含む。

【0003】 独自のデータの通信のためのATMネットワークを用いることができるアプリケーションは極めて多岐にわたる。ATMネットワークを用いることができるアプリケーションの大半はそのデータ要素の独自のフ

## 3

フォーマットを備える。これらは例えばインターネット・プロトコルのIPフォーマット・フレームまたはMPEG (moving picture export group) フォーマットを用いるフレームである。これらのアプリケーション・フレームのフォーマットとATMセルのフォーマットの間の適合はATMアダプテーション・レイヤすなわちAALと呼ばれるレイヤで実行される。より詳細に言えば、このレイヤはフレームをセルに分割し、その逆にネットワークから受信したセルをフレームに再組み立てする処理を担当する。

【0004】「実時間フロー」として知られる伝送されるデータ・フローの一部はネットワークがセルを提示する送信時間とジッタが最小になることを要求する。このケースは例えば電話データに関連する。電子メールなどのその他のデータ・フロー、以下、非実時間フローにはこれらの制約がない。実時間フローはネットワーク内の一定程度の優先度を利用する必要があり、このためにフローはリソースの予約による輻輳予防制御機構に支配される。

【0005】「サービス・カテゴリ」として知られるいくつかの主要フロー・クラスが、セル損失率、転送時間、ジッタ、最小伝送速度などのサービス品質パラメータが定義する、ユーザが所望のフローに関して追求するサービス品質に関してユーザが出す異なる要求を考慮し、ピーク伝送速度、平均伝送速度、バーストの最大サイズなどのトラフィック・パラメータが形成するこのフローの伝送速度の異なる特性を考慮する規格（[UIT-T, I. 371]、[AF-TM4. 0]）で定義されている。

【0006】実時間フローはCBR（固定伝送速度）またはVBRrt（実時間可変伝送速度）として知られる以下のカテゴリの1つに属する。非実時間フローはサービス品質がユーザ側のいかなる要求条件の目的でもないUBR（無指定伝送速度）、統計的特性が十分な精度で損失率に関する保証が可能なVBRnrt（非実時間可変伝送速度）または最小伝送速度もしくはネットワークの表示によるエンドツーエンドのフロー制御の交換条件としての低損失率が保証される利用可能伝送速度（ABR）として知られる以下のカテゴリの1つに属する。

【0007】すべてのATMスイッチは、図1aに示す方法で4つの主要な機能セット、ATMスイッチの各ポートへのアクセス機能1、ATMレイヤ機能2、スイッチ・ファブリック機能3および管理機能4を実施する。

【0008】アクセス機能1は前記ポートに接続された伝送媒体に適したフォーマットへのATMセルの変換およびその逆の処理に備える。この機能によって、伝送速度および着信セルを送信する伝送媒体が光、電気、無線、またはその他のタイプであるかどうかにかかわらず、着信セルを単一のフォーマットでATMレイヤへ向けることができる。スイッチのポートによっていくつか

## 4

のスイッチを接続できるがまたATMサービスのユーザをスイッチに接続することもできる。

【0009】アクセス機能において実施すべき処理動作はANSIとUITおよびATMフォーラム・システムに関する大量の標準設定の文献に記載されている。これらのドキュメントに記載されたインタフェースの主要クラスは次の通りである。

ドキュメントUIT-T G. 804、G. 703に規定されたPDH（独立同期デジタル・ハイアラキ）

## 10 インタフェース

ドキュメントUIT-T G. 708等に規定されたSDH（同期デジタル・ハイアラキ）インタフェース  
ドキュメントANSI-T 1. 105等に規定されたSONET（光同期ネットワーク）インタフェース  
ドキュメントaf-phy-0040. 000に規定された25. 6Mビット/s IBMインタフェース

【0010】ATMレイヤ機能2はいくつかの機能、特にセル・ヘッダの管理、VPI/VCI（仮想経路識別および仮想チャネル識別子）論理チャネルの翻訳、OAM（動作、運用および保守）管理セルの処理、UPC（使用パラメータ制御）、SCD（選択的セル廃棄）、EPD（早期PDU廃棄）、RM（リソース管理）セルとして知られるサブ機能を含むトラフィック管理として知られるトラフィック管理の主要部分を組み合わせる。

【0011】ATMレイヤ機能において実施すべき処理動作はUITおよびATMフォーラムの以下の標準ドキュメントに記載されている。

— B-ISDN ATN Layer Specification [UIT-T I. 361]

30 — B-ISDN Operation and Maintenance Principles and Functions [UIT-T I. 610]

— Traffic Management Specification Version 4. 0 [AF-TM 4. 0]

【0012】スクランブル機能3は、論理チャネルの翻訳（translation）中にATMレイヤが作成する表示の機能としてセルを入力方向から1つまたは複数の出力方向へ切り替える。

40 【0013】この機能はあらゆるATMスイッチの中核にあり、本明細書で引用する必要がない大量の文献で扱われている。スクランブル・リングおよびスクランブル・ネットワークはこの機能の2つの頻出するタイプの実施形態を構成する。

【0014】管理機能4は、交換仮想回線などの設定に必要なネットワークの集中監視エンティティとのインタフェースをとるスイッチのローカル監視（アラーム、スイッチおよびローカル・トポロジの構成の検出、バージョン管理など）などのサブ機能を含む。

50 【0015】これらのサブ機能の詳細な説明について

は、ATMフォーラムの標準設定の文献を参照されたい。

— ATM User-Network Interface (UNI) Signalling Specification Version 4.0 (af-sig-0061.000)

— Private Network-Network Interface Specification Version 1.0 (af-pnni-0055.000)

— Integrated Layer Management Interface (af-ilmi-0065.000)

【0016】これらのさまざまな機能は以下に記述するように相互にインタフェースされている。管理機能は、ATMレイヤへの接続がスイッチの外部ポートを経由しないのでアクセス機能が不要である点を除き、ユーザと同じように振る舞う。これと対照的に、管理機能はATMセルだけを処理せず、したがって、追加機能、すなわち適合機能であるAAL(ATM適合レイヤ)を用いてセグメント化し再組み立てする必要があるメッセージをも処理する。

【0017】ATMスイッチは集中アーキテクチャまたはゆるやかな分散アーキテクチャを備えたスイッチであることが多い。すなわち、スイッチ自体の機能はマイクロプロセッサが形成する計算容量、メモリが形成する記憶容量およびスイッチ・ファブリック内でセルをルーティングする容量を組み合わせた単一のハードウェア要素によって実行される。しかしながら、この集中はスイッチのモジュール的な性質とその構成要素の1つが障害になった場合でも動作が停止しないという能力に悪影響を与える。

【0018】標準的な解決策によると、必要に応じて2重化して同じ性質の障害要素をバックアップすることができ別個のハードウェア要素間にこれらの機能が分散される。これらのハードウェア要素はこの地点のネットワーク構成の関数として予見可能な処理負荷に対処するのに十分な数だけスイッチ内に実装される。実際にはこれらの要素はトレイ内に組み立てられトレイの底に敷設された1つまたは複数のデータ・バスによって相互にインタフェースをとる電子コンポーネントの実装基板である。それらは一般に「分散アーキテクチャ」と呼ばれる存在を画定する。

【0019】従来、図1bに示すような分散ATMスイッチのハードウェア・アーキテクチャでは3つのタイプのモジュールを区別している。それらはスイッチ・ファブリック・モジュール5、管理モジュール6、および回線カード・モジュール7<sub>1</sub>...7<sub>n</sub>である。スイッチの機能はこれらのさまざまなモジュールを介して分散されるが、回線カード・モジュールが少なくともアクセス

機能を扱い、スイッチ・ファブリック・モジュール5がスクランブル機能を扱い、管理モジュール6が管理機能を扱うという制約がある。

【0020】図1bでは、各回線カード・モジュールとスクランブル・モジュール間に存在するリンク

8<sub>1</sub>...8<sub>n</sub>は「スイッチ・ファブリック・ジャンクション」と呼ばれる。さらに、それぞれの回線カード・モジュールは1つまたは複数のポートを管理する能力があるアクセス機能を実施する。セルがスイッチを通過する際、まずセルの入力回線カードと呼ばれる第1の回線カードを通過し、次に出力回線カードと呼ばれる第2の回線カードを通過する。いくつかの入力回線カードは1つの同じ出力回線カードに向けて同時にセルを送信できるため、この出力回線カードの限られた出力伝送速度が原因でセル内部に輻輳が発生する可能性がある。セルを記憶して待ち行列に入れる機構が次に起動され、輻輳の解除はペンディングにされる。これらの記憶機構はスイッチ・ファブリックまたはこれらの要素のいくつかに同時に存在する入力部または出力部に見つけられる。したがって、使用される用語は「入力部における記憶」、「出力部における記憶」を備えたアーキテクチャなどである。

【0021】通信ネットワークのユーザはそのデータ要素のいくつかの交換モードを考慮することができる。これらのデータ要素を図2a~図2fに示す。図2aに示す二地点間モードは2人のユーザA、Dだけをリンクする。各ユーザは送信側および受信側である。このモードでは、ユーザの一方が送信したデータは他方のユーザによって受信される。二地点間モードの1つの変形形態は2人のユーザそれぞれの送信側と受信側の役割を固定することである(片方向二地点間通信)。

【0022】一地点多地点間モード(図2b)は1人のユーザが送信側専用で他のユーザが受信側専用である3人以上のユーザA、C、Dをリンクする。送信側が送信したデータは他のすべての受信側によって受信される。

【0023】多地点一地点間モード(図2c)も1人のユーザが受信側専用で他のユーザがすべて送信側専用である3人以上のユーザA、B、Cをリンクする。送信側の1人が送信したデータは受信側によって受信される。

【0024】最後に、多地点間モード(図2d)はそれぞれが送信側と受信側になれる少なくとも2人のユーザA、B、C、Dをリンクする。この最後のモードでは、ユーザのいずれかの1人が送信したデータは他のすべてのユーザと送信側によって受信される。

【0025】多地点間通信および一地点多地点間通信は図2eに示すイーサネット・ネットワークなどの共有媒体通信ネットワークの場合に特に適している。実際、この場合、すべてのユーザは単一の媒体に接続され、この媒体に接続されたすべての接続局A、B、C、Dは他のすべての局から送信されるメッセージをすべて受信す

る。これと対照的に、図2fに示すATMネットワークの場合は、1人のユーザが送信するセルのいくつかの宛先A、B、C、Dへの分散はネットワーク自体が問題のセルのコピーを生成することを要求する。

【0026】「接続」という用語は十分に定義されたユーザのセット間の上記のモードの1つによる、サービス品質パラメータ、トラフィック・パラメータなどの属性の特定のリストを授与されたあらゆる通信に適用される。

【0027】ATMネットワーク内での上記の異なるモードにおける通信の実施は、具体的にはデータ要素の信号方式、ルーティング、搬送およびリソースの管理といったいくつかの観点から考察できる。

【0028】二地点間接続に関しては、信号方式およびルーティングの態様が標準設定文庫のドキュメント

(UIT-T Q. 2931)、[AF-SIG 4.0]、[AF-PNNI 1.0]、[AF-IISP])に詳細に記載されている。

【0029】これらの態様はネットワーク内での接続のサービス品質およびトラフィックの制約を満足する2人のユーザ間の経路の決定からなる。経路は幹線のリストによって特徴付けられる。経路のそれぞれのスイッチは接続のスイッチへの入力幹線に関する論理チャネル番号を接続に割り当て、この識別子についてセルがたどる送信方向と次のスイッチでの接続の論理チャネル識別子の対応を得る翻訳テーブルを維持する。したがって、接続のすべてのセルはセル・ヘッダおよびローカル翻訳テーブルにある論理チャネル識別子を問い合わせることによってのみ一地点から次の地点へルーティングすることができる。

【0030】図1bに示すような分散アーキテクチャ・スイッチでは、この翻訳は入力回線カードのATMレイヤ機能によって実行される。このセルは次にスイッチ・ファブリックがセルを交換する先の出力スイッチ・ファブリック・ジャンクション表示を備えたスイッチ・ファブリック・モジュールへ渡される。この表示はセルの先頭に付加された特定のヘッダによって搬送される。この例による翻訳装置はフランス特許出願第2 670 972、2 681 164、2 726 669、および未公告の特許出願FR 97 07355に記載されている。

【0031】リソース管理の態様に関しては、接続のサービス・カテゴリによっていくつかのケースが可能である。

【0032】実時間フロー、すなわちCBRおよびVBR接続はリソースの予防的予約の対象である。したがって、実時間フローに起因する出力回線カードの輻輳の確率は低い。その結果、入力部に記憶メモリを構想することは一般に不要である。反対に、小型の出力記憶メモリがいくつかのエントリから同時に着信するセルを吸

収するために必要になる。

【0033】サービス品質が保証される非実時間フロー、すなわちVBR nrtまたはABR接続はある種の予防的予約処置の対象となり得るが、これらの処置はVBR nrtソースの極めて散在的な性質とABRソースが利用可能な帯域幅をすべて占有する傾向があるとすれば不十分である。したがって、このタイプのフローがVBR nrt接続の統計的な性質によって、またはABR接続のエンドツーエンドのフロー制御によってソースの総伝送速度が平均して利用可能伝送速度より小さいかそれに等しい状態を保つことが保証されるならば、これらのフローに関して、一時的な記憶メモリを準備する必要がある。

【0034】サービス品質の保証がない非実時間フロー、すなわちUBR接続に関しては、予防機構は不可能である。したがって、記憶メモリを準備してこの空間を異なる接続で共用して1つの接続が利用可能な空間を過度に占有することを防止する(公平の観点から)ことが必要である。また、UBR接続の処理が保証されたサービス品質を備えたフローの実行を阻害するおそれがある状況を防止する必要がある。

【0035】最後に挙げたポイントは、既存のATMスイッチの大半がすべてのタイプのフローに分けてリソースのプールを設定しているという事実によってそのすべての重要性を獲得する。これは特にスクランブル容量にあてはまる。しかしながら、これは記憶容量にもあてはまる、というのは利用可能なメモリをサービス・カテゴリの機能としてセグメント化するのは有利でないからである。したがって、フロー制御機構を実施してダウンラインのあらゆる輻輳を防止することを可能にする(出力回線カードは入力回線カードに遡及規則を適用してその伝送速度を調整する)かどうかはスイッチ次第である。この機構は保証されたサービス品質を備えたフローを阻害してはならない。さらに、この機構はダウンラインの輻輳をアップラインへ移行する傾向があるため、アップライン記憶リソース割り当てのポリシーに違反するセルを拒絶するのはATMレイヤ機能のサブ機能であるトラフィック管理機能次第である。

【0036】フランス特許出願第2 740 283は本出願人に代わって図1のこの種のフロー制御機構の例を提示する。この基本機構には、出力ジャンクションの機能としての入力メモリ内の特定の待ち行列に基づく

「ヘッドオブライン・ブロッキング(head-of-line blocking)」を解消する補助機構が関連付けられることが多い。実際、これらの機構は保証されたサービス品質のフローにはほとんど不可欠であることがわかっている。このような機構のさまざまな例はHPC ASIA 97のM. HYOJEONG SONG「A simple and fast scheduler for input queued AT

10

20

30

40

50



M switches」に記載され参照されている。

【0037】一地点多地点間接続を記号的に送信側ユーザを表す「根」と受信側ユーザを表す「葉」を備えた「木」で表すことができる。このタイプの接続の実施形態は信号方式およびルーティングに関して標準化されている。実行する必要があるのはまず二地点間接続を形成してからそれに新しい葉を付けることで一地点多地点間接続を形成することだけである。この葉の追加は根かそうでなければ葉がイニシアチブをとって実行できる。

【0038】セルの搬送に関して、二地点間接続モデル、すなわち入力翻訳だけを常に適用できるとは限らない。図1bの要素に似た要素が同じ基準で識別される図3aは一地点多地点間接続を示している。この接続はポートP1でスイッチに入り、ポートP2、P4、P5、P7でスイッチから出る。このケースでは、すでに考察した入力翻訳がスイッチ・ファブリック5にこの接続のそれぞれのセルを関連する3つのスイッチ・ファブリック・ジャンクション(73、74、7n)にコピーするよう命令できるが、セルの送信先の出力ポートを示すことはできない。これを行うには、この情報を翻訳テーブルに追加してそれを入力部から出力部まで送信する必要がある。さらに、出力論理チャネルは各出力ポートによって変わり、すべての出力部に単一の論理チャネルを割り当てることはできない。

【0039】こうした理由から、翻訳を2方向から行うことが一般に好ましい。それは、セルの論理チャネルをスイッチ内の接続を表す「同報通信インデックス」で置き換える入力翻訳と、次に同報通信インデックスを必要なコピーを備えた「ポート、論理チャネル」の式のペアのリストで表現する第2の出力翻訳である。

【0040】非実時間フローのリソース管理に関してはこれと逆に二地点間接続で実施されるフローの制御の機構を単に拡張することは望めない。実際、ヘッドオブライン・ブロッキング(head-of-line blocking)機構は各出力に特有の待ち行列を利用して可能な限りの出力セット数と同じ数、すなわち、N個のスイッチ・ファブリック・ジャンクションの場合は $2^{N-1}$ の待ち行列の管理を命令することになる。

【0041】別の解決策は、例えば、入力回線カード内にこのセルが関連する出力回線カードと同じ数のそれぞれのセルのコピーを作成することである。この手法はスイッチ・ファブリックのN個の周期にN個の出力ジャンクションへセルを送信させる必要があるという許容できない欠点があるが、これに反してスイッチ・ファブリック技術を用いると一般に単一の周期でセルをN個の出力へセルを転送できる。

【0042】これとは対照的に、本明細書で参照するフランス特許出願第2 740 283は伝送データを分離してそれらを1つずつ入力回線カードからスイッチ・ファブリックへ通過させ、すべての出力ポートでシェー

パとして知られる分離装置が指示する伝送速度に従って予約レベルが固定した固定値でリソースを完全に予約することで、異なる一地点多地点間UBR/ABRフローを1つの一地点多地点間CBRフローへローカルに変換する処理に基づく1つの解決策を記載する。

【0043】しかしながら、出力ポートの伝送速度が互いに非常に異なる場合にはこの手法には欠点がある。これはATMスイッチ内に25Mビット/秒のポートと622Mビット/秒のポートが混在するケースを考えると容易に分かる。実際このケースでは、シェーパの伝送速度を最も伝送速度が小さいポートの伝送速度に基づいて調整することが必要になろう。すべてのABR/UBR接続が同じシェーパを用いるため、すべての接続が阻害され、高速の伝送速度の出力ポートだけを使う接続さえも阻害される結果となる。

【0044】多地点一地点間接続および多地点間接続はATMを扱う標準化ユニットでは現在処理されていない。したがって、当座は、このタイプの接続について定義された信号方式またはルーティングは存在しない。

【0045】セルの搬送に関して、地理的発生地点が異なるデータ要素を1つの同じリンクに収束させる通信のあらゆるトポロジはいわゆる「PDU(プロトコル・データ単位)アプリケーション・フレームのインタレース」と呼ばれる問題を引き起こす。実際、PDUアプリケーション・フレームはAALレイヤによってセルにセグメント化されるため、異なるフレームのセルがインタレースされた形式で宛先へ送信される。フレーを再組み立てするには、宛先はそれぞれのセルが属するフレームを判定する能力を必要とする。AAL5適合レイヤで実施されるUBR接続で最も一般的に用いられるセグメント化機構ではこの識別が不可能である。この機構ではPDUフレームの最後のセルの識別だけが可能であるが、二地点間または一地点多地点間モードでは、ATMセルが順に送られるので、これで十分である。

【0046】集中アーキテクチャまたは出力部に記憶を備えたスイッチの場合、同じ1つのPDUフレームのセルをメモリに保存し、PDUフレームを完全に受信したらそれらをすべて送信することからなる手法で満足することが可能である。これは入力記憶を備えたスイッチには当てはまらない。というのは異なる回線カードからの伝送が調整されない場合にその異なる入力回線カードが共用するスイッチ・ファブリックへの同時アクセスは必ず別のインタレースを引き起こすからである。

【0047】上記の問題にもかかわらず、多地点一地点間および多地点間モードでの通信の必要が存在する。これらの問題は理論的には二地点間通信または一地点多地点間通信の重量によって処理できる。このケースは特にATMネットワーク内に共用媒体(ELAN)をエミュレートするために何か最新の機構が設置されるLANエミュレーションまたはLANE(ローカル・アレイ・ネ

ットワーク・エミュレーション)として知られるローカル・エリア・ネットワークのエミュレーション[AF LANE]に関するものである。参照できる例はLAN E規格で定義されるBUS(ブロードキャスト・オア・アンノウン・サーバ)として知られるサーバの例である。図3bにそのアーキテクチャを示すこのサーバを用いてユーザはエミュレートされた共用媒体のすべてのユーザへ、または直接接続されていない別のユーザへ向けてメッセージを同報通信することができる。図3bに示すように、ELANのそれぞれのユーザはBUSサーバに対して二地点間接続を確立し、BUSサーバはELANのすべてのユーザに対して一地点多地点間接続を確立する。従来技術では、BUSサーバ・タイプのブロードキャスト・サーバはATMネットワークのどこかに接続されたコストが高い専用ハードウェアまたはスイッチの管理モジュール上で実行される特異的プロセスによって実施される。最後の手法はELANのユーザの数が増加する場合には残念ながらあまり拡張できない。実際、新しいユーザがその相手に知られていない間、そのユーザが送受信するメッセージはBUSサーバを通過する必要がある。これは管理モジュールにとって致命的になり得る輻輳の原因であることが多い。

【0048】多地点一地点間通信および多地点間通信の必要の別の例をローカル・エリア・ネットワーク間のルーティングのエミュレーションによって説明する。この機能は特に異なるエミュレートされたローカル・ネットワーク(ELAN)間、またはELANの異なる仮想ローカル・エリア・ネットワーク(VLAN)[AFMPOA]間の仮想ルーティングを可能にするATMフォーラムのMPOA(マルチプロトコル・オーバーATM)として知られる規範に従って実施できる。ルーティングのエミュレーションを可能にする別の方法はATMスイッチの管理ユニットにルーティング・ソフトウェア・プログラムをインストールすることである。これに関して、仮想ルータはそれが相互接続する異なるELANのユーザに喩えることができる。これに基づいて仮想ルータは各ELANのLEC(LANエミュレーション・クライアント)として知られる機能を組み込む必要がある。仮想ルータはスイッチ内に例えば管理モジュールで実行される特異的なプロセスによって実施されなくてはならない。ELAN間の交換が十分に維持できるレベルに達した時にこの方法は前記モジュールに輻輳を引き起こす危険を伴う。

#### 【0049】

【発明が解決しようとする課題】本発明の目的は、上記の欠点を克服することである。

【0050】このため、本発明の目的は、分散アーキテクチャを備え入力部に記憶を備えたATMスイッチ内の、各入力回線カードへVBRnrt(非実時間可変伝送速度)カテゴリのフロー専用の合計平均伝送速度の関

数として調整された指定済みのn個のシェーパと、ABRおよびUBRカテゴリ・フロー専用で出力ポートの利用可能伝送速度(AvCR)の関数として調整されたn-1個の別のシェーパを分散することからなるフロー制御の方法である。

【0051】本発明のその他の特徴および有利な点は添付の図面に関する以下の説明から明らかになる。

#### 【0052】

【課題を解決するための手段】本発明による方法は、本明細書で上述した一地点多地点間接続のケースで参照される従来技術の欠点を克服する。その方法は特に図4に示す各入力回線カード7<sub>i</sub>上ではや1つのシェーパではなくいくつかのシェーパ9<sub>1</sub>~9<sub>n</sub>を実施することである。より詳細に言えば、1つのシェーパはカテゴリVBRnrtの実時間フロー専用である。このシェーパの伝送速度はVBRnrtフローの合計平均伝送速度の関数として調整される。UBR/ABRフロー専用の固定数n-1個のシェーパは、出力ポートの利用可能速度(AvCR)の関数として調整される。

【0053】ある動作範囲がその番号で識別される各UBR/ABRシェーパに割り当てられる。

【0054】動作範囲はシェーパが達成できる伝送速度の値を特徴づける。それらは伝送速度間隔である(例えば8~32Mビット/秒)。異なるシェーパ9<sub>1</sub>~9<sub>n</sub>の動作範囲は出力ポートの公称伝送速度の全範囲がカバーされるようにシステムの初期化の際に明確に割り当てられる。i+1ランクのシェーパの動作範囲はiランクのシェーパの範囲の下限よりも低い下限を備え、iランクのシェーパの範囲の下限よりもわずかに高い上限を備え、隣接するシェーパの動作範囲がわずかに重なるようになっている。

【0055】例えば、出力ポートの公称伝送速度が64kビット/秒から155Mビット/秒の間隔きざみになっているスイッチ内の次の分散を考えることができる。

番号	動作範囲
シェーパ1	155Mビット/秒~8Mビット/秒
シェーパ2	32Mビット/秒~2Mビット/秒
シェーパ3	8Mビット/秒~64kビット/秒

【0056】各一地点多地点間接続にはその入力回線カード(根)および出力ポート(葉)のデータ要素で構成される「木」が関連付けられる。(この用語は「根」と「葉」がユーザに適用される従来技術で用いる以前の用語と混同してはならない。)接続の生命が進行するにつれ、それに関連付けられた木は葉の追加または削除によって変更される。

【0057】本発明によれば、それぞれの木にはこの木に関連付けられたすべての一地点多地点間接続のセルを処理するシェーパが関連付けられる。シェーパは最も低い公称伝送速度を備えた「葉」のポートの公称伝送速度の関数として当業者の理解の範囲内の割り当て法則を用

いて選択される。上記の例に関して、以下のことが決定できる。

- すべての「葉」のポートが155Mビット/秒の公称伝送速度を備えた木をシェーパ1に割り当てる。
- 上記のカテゴリに入らず、「葉」ポートが25Mビット/秒未満の公称伝送速度を備えていない木をシェーパ2に割り当てる。
- 上記の2つのカテゴリに入らない木をシェーパ3に割り当てる。

【0058】このようにして、それぞれの一地点多地点間接続はシェーパに割り当てられ、この割り当ては接続の木の変更ごとに見直しができる。

【0059】シェーパ $9_1$ から $9_n$ への接続の割り当てが葉の公称伝送速度の関数として行われる一方、スイッチの伝送速度の調整は逆に葉の利用可能伝送速度の関数として行われる。利用可能伝送速度は公称伝送速度と予約伝送速度の差である。予約伝送速度の概念は接続CBR、VBR、ABR、またネットワークのポリシーがUBR接続の最小伝送速度の保証を意味する場合はUBRで補償される伝送速度の合計を表す実施態様に依存する伝送速度である。シェーパの伝送速度は関連するシェーパに割り当てられた木に属する葉の利用可能伝送速度の最小値をとって計算される。シェーパの伝送速度を調整するこの種の機構の正確な実施態様は当業者の理解の範囲内である。上記の例に関して、例えばスイッチの管理モジュールが各出力ポートの利用可能伝送速度を最新のものに保ち、それぞれの入力回線カードにそのシェーパに適用される伝送速度を周期的に送信することが想定できる。

【0060】VBR $n_{rt}$ 実時間フロー専用のシェーパについて言うと、これらは関連する接続の合計平均伝送速度の関数として調整される。

【0061】1つまたは複数の出力ポートでの利用可能伝送速度の大幅な低下に従って上記機構を適用することで、ランク $i$ シェーパに割り振られた伝送速度はその動作範囲と互換性がないことがある。したがって、この状況を回避するため、葉が関連するポートである木がそれによってより低い伝送速度に適合されたシェーパ、例えば $i+1$ ランクのシェーパに一時的に再割り当てされる例外的な機構を備えることが計画される。ここで動作範囲の重なりは自然なヒステリシス効果を生み、その結果、割り当ての変更は過度に頻繁には起こらず1つの同じ木に割り当てるシェーパの過度に頻繁な変更が防止できる。

【0062】異なる入力回線カード上にある同じ番号を備えたすべてのABR/UBRシェーパは1つの同じ出力ポートへそれぞれこのポートの利用可能伝送速度より低いかそれと等しい伝送速度でセルを同時に送信できる。この結果、最悪の場合に、入力回線カードからの着信伝送速度は関連するポートの利用可能伝送速度を容易

に超える可能性がある。この問題は異なる回線カードによるスクランブル・リソースへのアクセスの仲裁機構を実施することで解決される。この機構はまた多地点一地点間または多地点間接続に関するAAL5のセグメント化によるPDUアプリケーション・フレームのインターリーブに関する、参照される従来技術の欠点を克服することができる。

【0063】このために、ABR/UBR接続シェーパのそれぞれは、PDUアプリケーション・フレーム・モードで動作する優先待ち行列P1とセル・モードで動作する待ち行列P2の2つの待ち行列の間で仲裁を行う。待ち行列がPDUモードで動作すると言うことは、待ち行列に少なくとも1フレーム全体が収容され、その1つの同じフレームのセルが他のフレームのセルが挿入される可能性なしに順次送信されない限り、送信準備が整わないことを意味する。待ち行列がセル・モードで動作すると言うことは、待ち行列に少なくとも1つのセルが収容されると同時に送信準備が整うことを意味する。

【0064】次に、各入力回線カード上で、一地点多地点間接続で説明した方法と似た方法で多地点一地点間接続および多地点間接続が木に割り当てられる。多地点一地点間接続の場合、木は根だけになり葉は1枚である。多地点間接続の場合、木は考慮される入力回線カードを根として選択し、接続にかかわるすべてのポートを葉として選択できる。木は上記と同じ原理に従ってシェーパに割り当てられる。

【0065】多地点一地点間接続および多地点間接続は待ち行列P1だけを用いる。これはこの接続にはAAL5 PDUのインタレースの危険があるためである。

【0066】それぞれの回線カード $7_1$ のそれぞれのABR/UBRシェーパ $9_1 \sim 9_n$ はクラスに割り当てられる。クラスへの分散、すなわちすべてのシェーパの所与のランク $i$ への分割がある。この分割はシェーパのランクごと、すなわち動作範囲ごとに異なる。各クラスは次の原理で定義される。すなわち、所与の動作範囲で、この範囲のシェーパは、第1のシェーパに関連付けられた少なくとも1つの木と第2のシェーパに関連付けられた少なくとも1つの木が少なくとも1つの共通の葉を備えている場合、この範囲の別のシェーパと同じクラスに属する。逆に、2つのシェーパの1つに関連付けられた木のすべての葉が他方のシェーパに関連付けられた木のすべての葉との空白の交点を備える場合、図5aに示すように2つのシェーパは別のクラスに属する。この分割は動的であることに注意すべきである。すなわち、新しい木のシェーパへの割り当てまたはその解消には一般に分割の再計算が必要である。各分割は少なくとも1つのクラスを含み、最大でスイッチ内の回線カードと同じ数のクラスを含む。

【0067】上記の原理によるシェーパの分類のアルゴリズムは当業者の理解の範囲内である。

【0068】仲裁の原理は時間の $T = T1 + T2 + T3$ の長さの間隔への分割に基づく。同じクラスのさまざまなシェーパに《オーダ・オブ・スピーキング (order of speaking)》が割り当てられる。各時間間隔 $T$ で、《ライト・トゥ・スピーク (right to speak)》を備えた各クラスのシェーパが間隔 $T1$ にファイル $P1$ から抽出したPDUのセルの送信を開始できる。間隔 $T2$ で、シェーパは送信が開始されたが新しいPDUの送信を開始できないPDUのセルの送信を続行できる。この同じシェーパはもちろん $P1$ から受信したセルを送信できない( $P1$ の $P2$ に対する優先原理に従って)場合に間隔 $T1 + T2$ の間、 $P2$ から受信したセルを送信できる。その送信が $T1 + T2$ の終了前に終了すると、シェーパはライト・トゥ・スピーク (right to speak) をそれ自体も同じ間隔に送信が可能なそのクラスの別のシェーパへ引き渡すことができる。間隔 $T1 + T2$ が経過すると各回線カードは間隔 $T3$ 内で実行される収集機構によってメッセージを送信しようということを示す。間隔 $T$ の《ターン・トゥ・スピーク (turn to speak)》も表明された希望の関数としてこのPDUの終了前にPDUの送信を中断しなければならなかったシェーパのために計算される。《ターン・トゥ・スピーク (turn to speak)》を計算するアルゴリズムはスイッチの管理ポリシーが必要とする場合は回線カード間の一定の公平さを保証しなくてはならない。この種のアルゴリズムは当業者の理解の範囲内である。上記のことはすべてシェーパの所与のランク $i$ に関する。したがって、例えば間隔 $T$ 、 $T1$ 、 $T2$ 、 $T3$ は $i$ によって変化する、すなわち、考慮されるシェーパの動作範囲によって変化する。

#### 【0069】

【発明の実施の形態】トークンに基づく上記の原理の一実施態様を以下に例を用いて説明する。図5b、5cおよび5dに示すこの機構の異なる位相について説明する。

【0070】トークンは回線カード $7_i$ の1つかそうでなければ管理モジュール4かそうでなければスイッチ・ファブリック2自体であるマスタ・ユニットのイニシアチブによって異なる回線カード $7_i$ との間で交換される特殊セルを備える。本発明は3つのタイプのトークンを区別する。

【0071】図5bの《SYNC SIGNALトークン》は周期性 $T$ を備えたマスタ・ユニットによって一般同報通信で $i$ ランク・シェーパが少なくともそれに関連付けられた1つの木を所有するすべての回線カードへ周期的に送信される。このトークンは異なる回線カード間の送信信号の概算の同期化に用いられる。またこのトークンは各クラスの《ターン・トゥ・スピーク (turn to speak)》を搬送する。

【0072】図5cの《SPEECHトークン》は上記に識別した各クラスに特有である。このトークンはランク $i$ のシェーパが考慮されるクラスの《ターン・トゥ・スピーク (turn to speak)》の間隔 $T1 + T2$ の終了前に次の回線カードへの送信を完了した回線カードの1つによって送信される。

【0073】図5dの《COLLECTIONトークン》は周期 $T$ の最後にSYNC SIGNALトークンを受信した回線カードの1つによって送信される。このトークンは後続の周期に送信をするか送信がすでに開始しているPDUの最後まで送信を続行する必要がある標識を内部に記録するすべての関連する回線カード内で1つのユニットから次のユニットへと中継される。最終的に、トークンはマスタ・ユニットによって受信され処理される。

【0074】ブロードキャスト・グループの各ユーザのサーバへの二地点間接続と同じグループのすべてのユーザへのサーバの一地点多地点間接続の重畳に起因するBUSサーバなどのブロードキャスト・サーバのケースに関する従来技術の欠点を克服するため、ATMスイッチ内の異なる接続間に短絡が形成される。これによってPDUアプリケーション・フレームの中継がジャンクション装置によって完全に実行され、サーバはブロードキャスト・グループの管理、すなわちグループへの新しいユーザの到着またはグループからのユーザの出発に必要な動作以外の他のすべての役割から解放される多地点間通信を確立することが可能になる。回線カードのレベルでのPDUフレームの中継プロセスは上記の仲裁の実施とフロー制御プロセスによって可能になる。

【0075】図6aに参照10を持つブロードキャスト・サーバはスイッチの管理モジュール4上で実行されるソフトウェア・プロセスである。従来技術の動作モードでは、通信モジュールは1つのブロードキャスト・グループに属する3人のユーザの例を示す図6bに示されている。本発明による動作モードでは、通信モジュールは二地点間接続 $A1$ 、 $A2$ 、 $A3$ および一地点多地点間接続 $M$ がグループの各ユーザからすべての他のユーザへの一地点多地点間接続 $B1$ 、 $B2$ 、 $B3$ に置き換えられる図6cに示す通信モジュールである。

【0076】通信モジュールの変更はATMレイヤで実施され、図6aの回線カード $7_1$ 、 $7_2$ 、 $7_3$ のTRANS IN、TRANS OUTと参照される入力および出力翻訳の簡素な変換によって可能になる。 $A_i$ タイプの二地点間接続の入力翻訳について言えば、それらは従来技術にある。 $VLA_i \rightarrow (VLA'_i, LG)$ において $VLA_i$ は入力幹線の $A_i$ に関連付けられた論理経路識別子、 $VLA'_i$ は管理モジュール内のこの同じ接続に関連付けられた識別子、 $LG$ は管理モジュールのスイッチ・ファブリック・ジャンクションの識別である。これはすべてのPDUフレームが実施モジュール内のサ

サーバによって処理されることを表す。これと逆に、サーバから所与のブロードキャスト・グループのユーザへ再送されるすべてのフレームは一地点多地点間接続Mに特有の論理チャネルVLMで送信される。出力翻訳によって関連する各回線カード内で、それぞれのセルがいくつかのポートP1、P2、...、Ppから出力される必要がある場合にはそれぞれのセルのコピーが可能な方法で同報通信インデックスVLMを出力幹線上のこの接続に特有の論理チャネル識別子VLMiへ変換することが可能になる。したがって、最後の翻訳はVLM→(VLM P1、P1)、(VLM P2、P2)、...、(VLM Pp、Pp)の形式を備える。

【0077】本発明によれば、短絡によって変更される唯一の翻訳はVLAi→(VLM、L1、L2、L3、...、Lq)になる入力翻訳である。ただしL1からLqは接続Mに関連するすべてのスイッチ・ファブリック・ジャンクションを指定する。

【0078】本発明による動作モードでは、ブロードキャスト・サーバは従来技術と同じ方法でユーザの到着と出発を管理し、上記の通信モジュールに到着するために変更すべき入力翻訳を決定する。

【0079】図6dの表は図6aに示す4人のユーザui、u21、u22およびu3を含むブロードキャスト・グループのケースで短絡ありと短絡なしの翻訳例である。

【0080】注1：図6bは詳細な図ではないが接続A1、A2、A3がLANEサーバBUSには当てはまらない片方向接続のケースを表す。

【0081】注2：接続A1、A2、A3およびMは少なくとも信号方式に関して存在を継続する。

【0082】注3：アプリケーションによっては、サーバ自体がユーザ・グループの一部を形成する、または形成しないと考えることが常に可能である。

【0083】最後に、上記の仲裁機構に加えて、フレーム・ルータIPが管理モジュール内で実施された場合の上記に参照される従来技術の欠点を克服できる第3の機構を備えることが可能である。この方法を用いてルータの役割を経路の計算に限定することで(従来技術)、ATMスイッチ内で中継機能IPの実際の非集中化を実現できる。

【0084】図3cにこの方法が使用できるケースを示す。内部ルータはそれが認識できるELANと同じ数のLECを備える。ELANAに属するユーザuAがフレームIPをユーザuBへ送信しようとするれば、まずエミュレータELAN上の同報通信手段(BUSブロードキャスト・サーバ)を使用する。スイッチ内部のルータ機能がユーザuBの存在に気づいていれば、ELANAエミュレーションに関連付けられたこのルータの機能LEC Aはそれ自体をユーザuBへ向けられたすべてのIPフレームの宛先であると宣言する。その後、ユー

ザuAは直接のATM接続を用いてユーザuBへのフレーム送信のLANE規格に規定された一般手順に従ってLEC A機能を実行する。従来技術では、フレームをLECB機能とユーザuBの間に存在する直接のATM接続を用いてユーザuBへ向けて中継するにはフレームを管理モジュールにある内部ルータまで返送する必要がある。

【0085】本発明による方法は、LEC A機能に関する直接接続に到着するあらゆるPDUアプリケーション・フレームに関して、入力回線カードはその最初のセルを検査してそこから宛先のIPアドレスを抽出する。次に入力回線カードは管理モジュールから受信するルーティング情報によって更新されるキャッシュ・テーブルを検索し、テーブル内のIPアドレスの前にベア(論理チャネル、送信方向)を見つける。論理チャネルはLECBとuB間の直接接続が用いるチャネルである。送信方向はこの直接接続に関連するスイッチ・ファブリック・ジャンクションの識別子である。入力回線カードがキャッシュ・テーブル内でIPアドレスを検索した結果それが見つからない場合、入力回線カードは管理モジュールにキャッシュ更新要求を送信する。次にテーブル内で見つけられた情報は関連するPDUのそれぞれのセルのATMヘッダの翻訳に用いられる。したがって、動的翻訳が得られ、各PDUの通過時に翻訳テーブルが変更される可能性がある。結果として、「動的短絡」が2つの二地点間接続間に設定される。

【図面の簡単な説明】

【図1a】本発明によるATMスイッチの原理を示す図である。

【図1b】従来技術による分散アーキテクチャを備えたATMスイッチの原理を示す図である。

【図2a】ATMネットワークのユーザ間の通信モードを示す図である。

【図2b】ATMネットワークのユーザ間の通信モードを示す図である。

【図2c】ATMネットワークのユーザ間の通信モードを示す図である。

【図2d】ATMネットワークのユーザ間の通信モードを示す図である。

【図2e】ATMネットワークのユーザ間の通信モードを示す図である。

【図2f】ATMネットワークのユーザ間の通信モードを示す図である。

【図3a】一地点多地点間接続時のスイッチ内のATMセルのルーティング例を示す図である。

【図3b】エミュレートLANアーキテクチャでの二地点間接続または一地点多地点間接続の重畳の例を示す図である。

【図3c】ELAN間のルーティングの原理を示す図である。

【図 3 d】本発明が実施する短絡の原理を示す図である。

【図 4】本発明による木のシェーパへの割り当ての原理を示す図である。

【図 5 a】本発明によるシェーパの分類動作の例を示す図である。

【図 5 b】本発明によるトークンの同報通信の例を示す図である。

【図 5 c】本発明によるトークンの同報通信の例を示す図である。

【図 5 d】本発明によるトークンの同報通信の例を示す図である。

【図 6 a】本発明によるブロードキャスト・サーバの実施例を示す図である。

【図 6 b】本発明によるブロードキャスト・サーバの実

施例を示す図である。

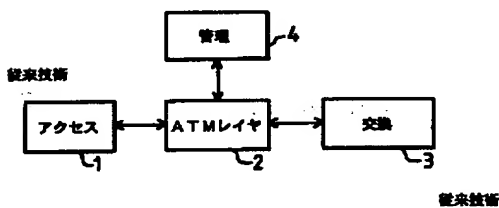
【図 6 c】本発明によるブロードキャスト・サーバの実施例を示す図である。

【図 6 d】本発明によるブロードキャスト・サーバの実施例を示す図である。

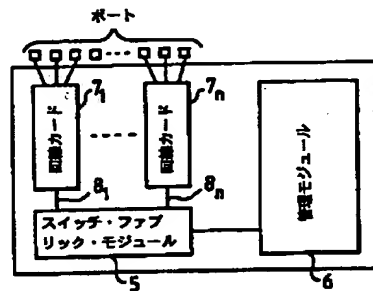
#### 【符号の説明】

- 1 ATMスイッチの各ポートへのアクセス機能
- 2 ATMレイヤ機能
- 3 スイッチ・ファブリック機能
- 10 4 管理機能
- 5 スイッチ・ファブリック・モジュール
- 6 管理モジュール
- 7 1... 7<sub>n</sub> 回線カード・モジュール
- 8 1... 8<sub>n</sub> リンク
- 9 1... 9<sub>n</sub> シェーパ

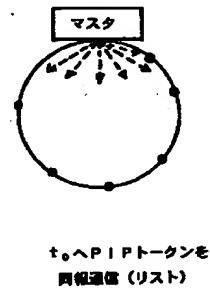
【図 1 a】



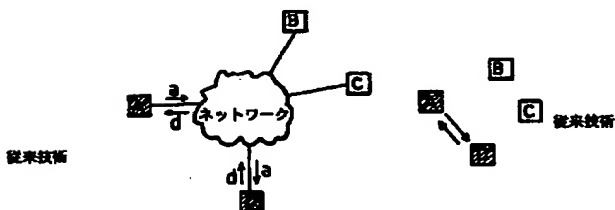
【図 1 b】



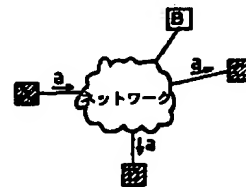
【図 5 b】



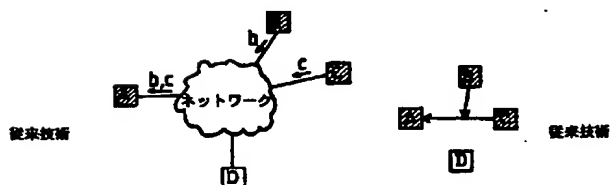
【図 2 a】



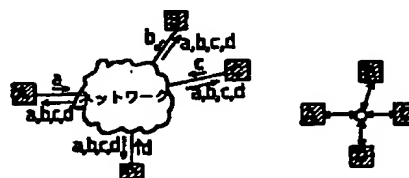
【図 2 b】



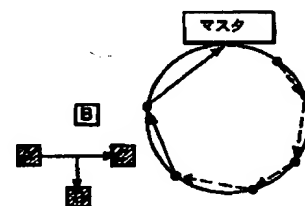
【図 2 c】



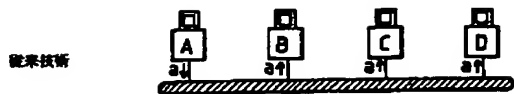
【図 2 d】



【図 5 d】

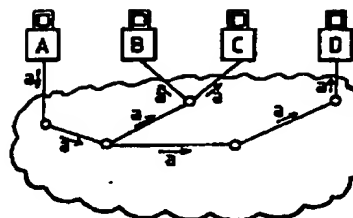


【図 2 e】



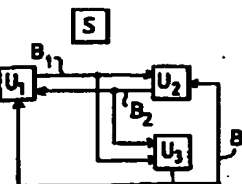
従来技術

【図 2 f】

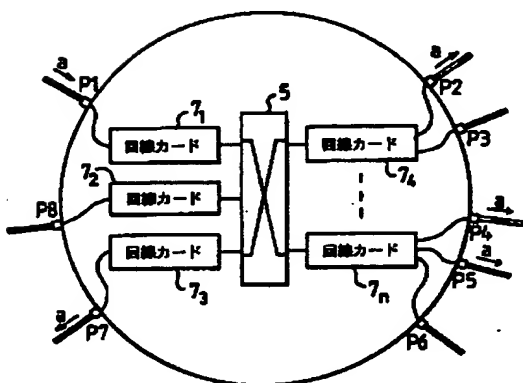


従来技術

【図 6 c】

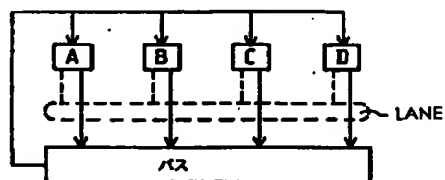


【図 3 a】

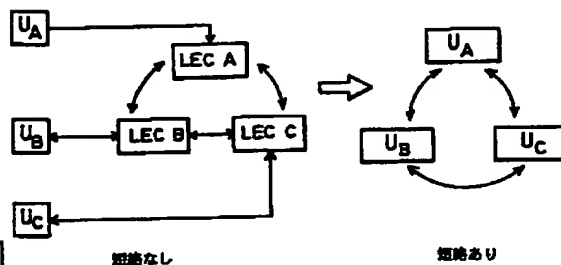


従来技術

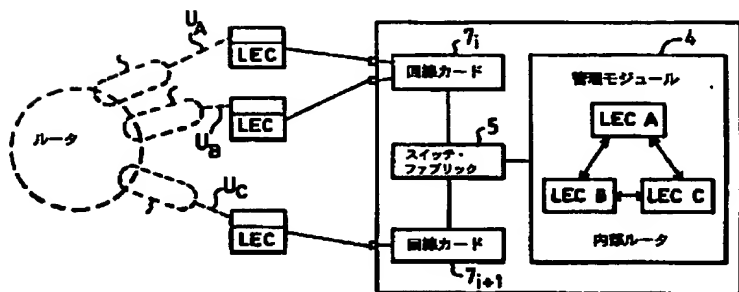
【図 3 b】



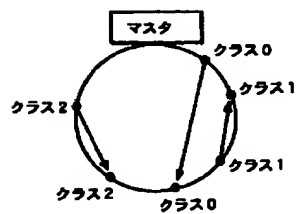
【図 3 d】



【図 3 c】

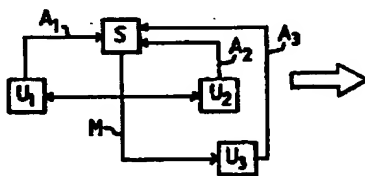


【図 5 c】

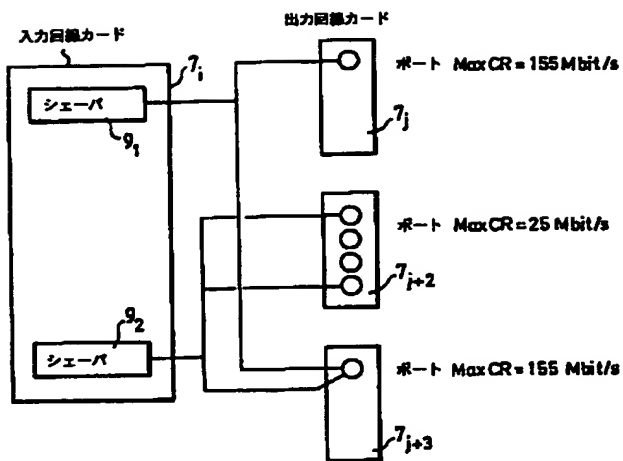


クラス別 SPEECH トークンを  
 $t < t_0 + T_1 + T_2$  へ送信

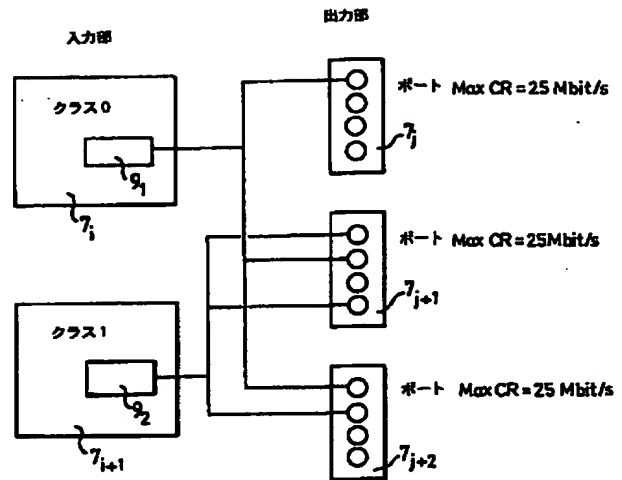
【図 6 b】



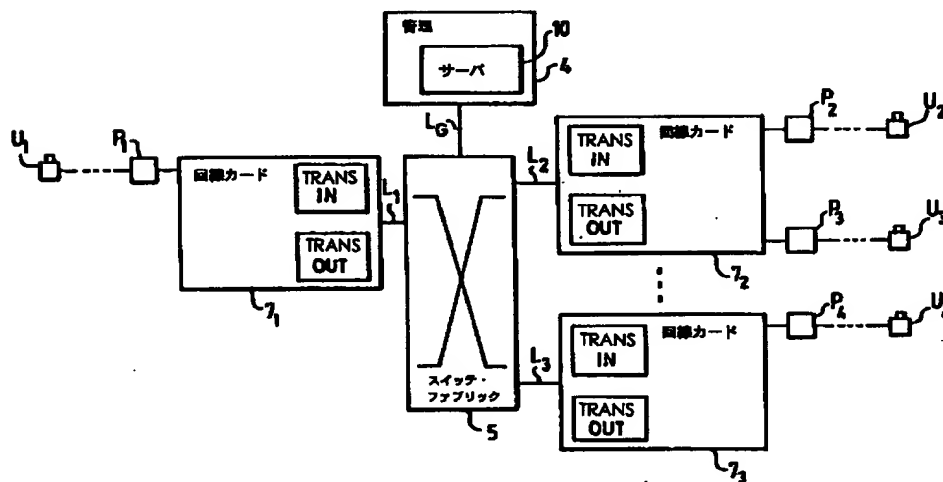
【図4】



【図5a】



【図6a】



【図6d】

	短絡なし	短絡あり
回線カード1、入力	$VLA\ 1 \rightarrow VLA'\ 1, LG$	$VLA\ 1 \rightarrow VLM, L1, L2, L3$
回線カード1、出力	$VLM \rightarrow VLM1, P1$	同上
回線カード2、入力	$VLA\ 21 \rightarrow VLA'\ 21, LG$ $VLA\ 22 \rightarrow VLA'\ 22, LG$	$VLA\ 21 \rightarrow VLM, L1, L2, L3$ $VLA\ 22 \rightarrow VLM, L1, L2, L3$
回線カード2、出力	$VLM \rightarrow (VLM\ 21, P21)$ $(VLM\ 22, P22)$	同上
回線カード3、入力	$VLA\ 3 \rightarrow VLA'\ 3, LG$	$VLA\ 3 \rightarrow VLM, L1, L2, L3$
回線カード3、出力	$VLM \rightarrow VLM3, P3$	同上



フロントページの続き

(72)発明者 ディディエ ゲラン  
フランス国, 78890 ガランスイエール,  
リュ サン ミシェル, 5 ビス番地

(72)発明者 フィリップ エロー  
フランス国, 92240 マラコフ, リュ  
ガリエニ, 44番地

## 【外国語明細書】

## 1 Title of Invention

METHOD FOR THE CONTROL OF FLOWS WITHIN AN ATM SWITCH WITH  
DISTRIBUTED ARCHITECTURE

## 2 Claims

1. A method for the control of flows within an ATM switch with distributed architecture and storage at input, wherein the method consists of the distribution, to each input line card, of a specified number  $n$  of shapers, a shaper being dedicated to the VBRnrt (variable bit rate non-real time) category flows and being adjusted as a function of the totalized mean bit rate and the  $n-1$  other shapers being dedicated to the ABR and UBR category flows, wherein the method further consists in associating, with each ABR or UBR connection, a tree whose root is constituted by its input line card and whose leaves are constituted by the output ports, then in assigning this tree to a shaper chosen as a function of the nominal bit rate of its leaves, the shaper being adjusted as a function of the available bit rate of the leaves of all the trees assigned to this shaper.

2. A method according to claim 1, consisting in assigning, to each shaper, a first queue  $P_1$  working in PDU application frame mode and a second queue  $P_2$  working in cell mode.

3. A method according to claim 2, consisting in reserving the first queue  $P_1$  to the multipoint-to-point and multipoint-to-multipoint connections.

4. A method according to one of the claims 1 to 3, consisting in assigning each shaper to a class defined by its range of operation, a shaper of this range being in the same class as another shaper if the respective trees associated with each of the shapers has at least one common leaf

5. A method according to claim 4, consisting in assigning an « order of speaking » to the shapers of one and the same class, in limiting the right to speak to within a specified time interval  $T$ , divided into three elementary time intervals  $T_1$ ,  $T_2$  et  $T_3$  such that  $T = T_1 + T_2 + T_3$ , the beginning of a transmission starting in the interval  $T_1$ , the end of a transmission taking place in the interval  $T_1 + T_2$ , the turn to speak in the next interval  $T$  being determined as a function of wishes expressed by each shaper during the interval  $T_3$  in favoring the shaper that has suspended the transmission of a PDU application frame.

6. A method according to claim 5, consisting in designating a « master » line card for the transmission of a « SYNC SIGNAL token » with a periodicity T, in a general broadcast to all the line cards whose rank i shaper possesses at least one tree that is associated with it, bringing about the transmission of a « speech » token by one of the line cards whose rank i shaper has finished sending before the end of the interval T1+T2 to the next line card with the turn to speak in the class considered and in bringing about the transmission, by the line card that has received the SYNC SIGNAL token at the end of the period T during the time interval T<sub>3</sub>, of a « collection » token towards the master line card in which there is recorded an indicator indicating its need to transmit during the following time interval T.

7. A method according to claim 6, consisting in carrying out communications in multipoint-to-multipoint mode by shifting the frame relaying function of the broadcasting server in the ATM layer of the line cards by modification of the input translations of the point-to-point connections (Ai), which go from  $VLAi \rightarrow (VLA'i, LG)$  to  $VLAi \rightarrow (VLM, L1, L2, L3, \dots, Lq)$  where VLAi designates the logic channel identifier associated with Ai on the input artery, VLA'i designates the logic channel identifier associated with this same connection in the management module, LG designates the identity of the switch fabric junction of the management module, VLM designates the broadcasting index of the point-to-multipoint connection M and L1, L2, ..., Lq designate the identity of the switch fabric functions concerned by the connection M.

8. A method according to claim 6, consisting in carrying out communications in multipoint-to-multipoint mode by shifting the frame relaying function of a routing emulation device in the ATM layer of the line cards by an examination of the first cell of each PDU frame in order to extract the IP address of the addressee therefrom, by making a search in a cache table of a pair formed by the logic channel and the outgoing direction before the IP address concerned, and using the translation obtained in all the cells of the PDU frame, the cache table being updated by means of routing information coming from the routing emulation function that resides in the management module and in transmitting a request for the updating of the cache to the management module if the desired IP address is not therein.

### 3 Detailed Description of Invention

#### BACKGROUND OF THE INVENTION

##### 1. Field of the Invention

The present invention relates to a method for the control of flows for point-to-point, point-to-multipoint, multipoint-to-point and multipoint-to-multipoint non-real time connections within an ATM switch with distributed architecture and storage at input.

The communication networks known as ATM or asynchronous transfer mode networks enable the circulation of fixed length packets known as ATM cells consisting of a 5-byte header and a 48-byte body. The header contains in particular a logic channel identifier known as a VPI/VCI (or virtual path identification and virtual channel identifier) field that enables the routing of the cells in the switches that it encounters on its path between the sender user and the addressee user.

The applications capable of using ATM networks for the communication of their data are highly varied. Most of the applications capable of using ATM networks have their own format for their data elements: these may be for example IP format frames of the Internet protocol or else frames using the format of the MPEG (moving picture export group) format. The adaptation between the format of these application frames and the format of the ATM cells is done in a layer called an ATM adaptation layer or AAL. In particular, this layer is responsible for segmenting the frames into cells and conversely reassembling the cells received from the network into frames.

Some of the data flows conveyed, known as "real time flows" require that the transit time and the jitter to which the network subjects their cells should be the minimum. This case relates for example to telephone data. Other data flows, hereinafter called non-real time flows, for example electronic mail, do not have these constraints. The real time flows must benefit from a certain degree of priority within the network and, for this purpose, they are subjected to a preventive congestion control mechanism by means of the reservation of resources.

Several major classes of flows known as "service categories" are defined in the standards ([UIT-T.I.371], [AF TM4.0]) to take account of the

different demands that users may make as regards the quality of the service that they seek for a desired flow, defined by service quality parameters such as: cell loss rate, transfer time, jitter, minimum bit rate, etc. and to take account of the different characteristics of the bit rate of this flow, formed by traffic parameters such as peak bit rate, mean bit rate, maximum size of a burst, etc.

The real-time flow comes within one of the categories known as CBR (constant bit rate) or VBRrt (variable bit rate - real time). The non-real-time flows come within one of the following categories known as UBR (unspecified bit rate) whose service quality is not the object of any requirement on the part of the user, VBRnrt (variable bit rate - non-real-time) whose statistical characteristics are known with sufficient precision to enable a guarantee with respect to the loss rate or available bit rate (ABR) for which a minimum bit rate is guaranteed or a low loss rate as a trade-off for an end-to-end flow control according to the indications of the network.

Any ATM switch, in the manner shown in Figure 1a, implements four major sets of functions, an access function 1 to each port of an ATM switch, an ATM layer function 2, a switch fabric function 3 and a management function 4.

The access function 1 provides for the conversion of the ATM cells into the format that is suited to the transmission medium connected to said port and vice versa. This function makes it possible to present incoming cells to the ATM layer in a single format that is independent of the bit rate and the optical, electrical, radio or other type of technology of the transmission medium from which they come. The ports of a switch enable the connection of several switches together but they also enable the connection of a user of ATM services to a switch.

The processing operations to be implemented in the access function are described in a huge volume of standard-setting literature on the ANSI as well as the UIT and ATM Forum systems. The major classes of interface defined in these documents are:

The PDH (plesiochronous digital hierarchy) interface defined in the document UIT-T G.804, G.703.

The SDH (synchronous digital hierarchy) interface defined in the document UIT-T G.708, etc.

The SONET (synchronous optical network) interface defined in the document ANSI-T 1:105, etc.

The 25.6 Mbit/s IBM interface defined in the document af-phy-0040.000.

The ATM layer function 2 combines several functions especially the management of the cell headers, the translation of the VPI/VCI (virtual path identification and virtual channel identifier) logic channels, the processing of the OAM (operations, administration and maintenance) management cells, a major part of the management of the traffic known as traffic management comprising sub-functions known as UPC (usage parameter control), SCD (selective cell discard), EPD (early PDU discard), RM (resource management) cells, etc.

The processing operations to be implemented in the ATM layer function are described especially in the following standard documents of the UIT and the ATM Forum:

- B-ISDN ATN Layer Specification [UIT-T I.361]
- B-ISDN Operation and Maintenance Principles and Functions [UIT-T I.610]
- Traffic Management Specification Version 4.0 [AF-TM 4.0]

The scrambling function 3 switches the cells from an input direction to one or more output directions, as a function of indications prepared by the ATM layer during the translation of the logic channels.

This function is of the core of any ATM switch and has been dealt with in a large body of literature that need not be recalled here. The scrambling ring and the scrambling network constitute two frequent types of implementation of this function.

The management function 4 comprises sub-functions such as: the local supervision of the switch (alarms, discovery of the configuration of the switch and of the local topology, management of versions, etc.), interfacing with the centralized supervision entity of the network, the interfacing needed to set up switched virtual circuits, etc.

For a more detailed description of some of these sub-functions, reference will be made to the standard-setting literature of the ATM Forum:

- ATM User-Network Interface (UNI) Signalling Specification Version 4.0 (af-sig-0061.000)

- Private Network-Network Interface Specification Version 1.0 (af-pnni-0055.000)

- Integrated Layer Management Interface (af-ilmi-0065.000)

These different functions are mutually interfaced as indicated here below. It must be noted that the management function behaves exactly like a user except that its connection to the ATM layer does not go through an external port of the switch and therefore does not require any access function. By contrast, the management function does process ATM cells alone but also messages which it must therefore segment and reassemble by means of an AAL (ATM adaptation layer) which is therefore an additional function: the adaptation function.

The ATM switches are often switches with centralized architecture or weakly distributed architecture, that is to say the functions of the switch itself are performed by a single hardware element that combines computation capacities formed by microprocessors, storage capacities formed by memories and capacities for routing the cells in the switch fabric. However, this concentration adversely affects the modular nature of the switch and its ability to remain functional when one of its constituent elements goes out of order.

According to a standard solution, the functions are distributed among distinct hardware elements that may, if necessary, be doubled to enable the back-up of a faulty element of the same nature. These hardware elements installed in the switch in sufficient numbers to cope with the processing load that is foreseeable as a function of the configuration of the network at this place. In practice, these elements are electronic component boards assembled in a tray and interfacing with one another by means of one or more data buses placed at the bottom of the tray. They define what is common called a "distributed architecture".

Conventionally, the hardware architecture of a distributed ATM switch, as can be seen in Figure 1b, distinguishes between three types of modules: a switch fabric module 5, a management module 6 and line card modules  $7_1 \dots 7_n$ . The functions of the switch are distributed over these different modules with, however, the constraint wherein the line card modules deal at least with the access function, the switch fabric module 5

with the scrambling function and the management module 6 with the management function.

In Figure 1b, the links  $8_1 \dots 8_n$  existing between each line card module and the scrambling module are called "switch fabric junctions". Furthermore, each line card module implements an access function capable of managing one or more ports. When a cell goes through a switch, it begins by crossing a first line card called an input line card for the cell and then a second line card called an output line card. Since several input line cards may simultaneously send cells towards one and the same output line card, there may be congestion in the cells owing to the limited output bit rate of this output line card. Mechanisms for storing the cells and placing them in queues are then activated, pending the clearing of the congestion. These storage mechanisms may be found at the input or the output, in the switch fabric or in several of these elements at the same time. The term used is then architecture with "storage at input", "storage at output", etc.

The users of a communications network may envisage several modes of exchange of their data elements. These data elements are shown schematically in Figures 2a to 2f. The point-to-point mode shown in Figure 2a links up two users A, D exclusively. Each of them may be a transmitter and a receiver. In this mode, all that is sent by one of the users is received by the other user. One variant of the point-to-point mode consists in specializing the sender and receiver roles of each of the two users (one-way point-to-point communications).

The point-to-multipoint mode (Figure 2b) links up more than two users, A, C, D, one of whom is exclusively a sender while the others are exclusively receivers. All that has is sent by the sender is received by all the receivers.

The multipoint-to-point mode (Figure 2c) also links up more than two users A, B, C one of whom is exclusively a receiver while the others are exclusively senders. All that has sent by one of the senders is received by the receiver.

Finally, the multipoint-to-multipoint mode (Figure 2d) links up at least two users A, B, C, D each capable of being a sender and a receiver. In this last-named mode, all that has been sent by any one of the users is received by all the other users and also by the sender.



The multipoint-to-multipoint communications and the point-to-multipoint communications are especially natural in the case of a shared medium communications network such as the Ethernet network shown schematically in Figure 2e. Indeed, in this case, all the users are connected to a single medium and all the connected stations A, B, C, D connected to this medium receive all the messages sent by the other stations. On the contrary, in the case of an ATM network as shown schematically in Figure 2f, the distribution, to several addressees A, B, C, D, of a cell sent by one of the users requires that the network should itself generate copies of the cell in question.

The term 'connection' is applied to any communication according to one of the modes defined here above, between a well-defined set of users, this communication being endowed with a specific list of attributes such as: service quality parameters, traffic parameters, etc.

The implementation, within an ATM network, of communications in the different modes defined here above may be seen from several points of view, in particular signalling, routing, conveying of data elements and management of resources.

With regard to point-to-point connections, the signalling and routing aspects are extensively described in the documents ([UIT-T Q.2931], [AF-SIG 4.0], [AF-PNNI1.0], [AF-IISP]) of the standard-setting literature.

These aspects consist of the determining, in the network, of a path between the two users such that this path meets the constraints of service quality and traffic of the connection. The path is characterized by a list of arteries. Each switch of the path assigns the connection a logic channel number pertaining to the input artery of the connection into the switch and maintains a translation table which, for this identifier, obtains a correspondence between the outgoing direction that the cell must take and the logic channel identifier of the connection in the next switch. Thus, every cell of the connection may be routed from one point to the next solely by consulting the logic channel identifier present in the cell header and the local translation table.

In a distributed architecture switch, such as the one of Figure 1b, this translation can be done by the ATM layer function of the input line card. The cell is then handed over to the switch fabric module with an indication of

the output switch fabric junction towards which the switch fabric must switch the cell. This indication may be conveyed by a specific header added to the beginning of the cell. Translation devices in accordance with this example have been described by the Applicant, for example in the French patent applications Nos. 2 670 972, 2 681 164, 2 726 669, and the as yet unpublished patent application FR 97 07355.

With regard to the resource management aspect, several cases are possible depending on the service category of the connection.

The real-time flows, namely the CBR and VBRrt connections, are subjected to a preventive reservation of resources. The probability of congestion of an output line card because of real-time flows is therefore low. Thus, it is generally not necessary to plan for storage memories at input. By contrast, a small-sized output storage memory is necessary to absorb cells coming simultaneously from several entries.

The non-real-time flows with guarantee of service quality, namely the VBRnrt or ABR connections, can be subjected to certain measures of preventive reservation, but these measures are insufficient given the highly sporadic nature of the VBRnrt sources and the tendency of the ABR sources to occupy the entire available bandwidth. For these flows, it is therefore necessary to provide for temporary storage memories, given that this type of flow, by its statistical nature for the VBRnrt connections or by its end-to-end flow control for the ABR connections, ensures that the total bit rate of the sources remains on an average smaller than or equal to the available bit rate.

For the non-real time flows without guarantee of service quality, namely the UBR connections, no preventive mechanism is possible. It is therefore necessary to provide also for storage memories and share this space between the different connections in preventing one connection from taking up too much available space (in accordance with the notion of fairness). It is also necessary to prevent a situation where the processing of the UBR connections could hamper the running of flows with guaranteed service quality.

This last-named point acquires all its importance through the fact that most of the existing ATM switches set up a pooling of resources between all types of flows. This is especially true with respect to scrambling

capacities. However, it is also the case with storage capacities for it is not advantageous to segment the available memory as a function of the service category. It is therefore up to the switch to implement a flow control mechanism making it possible to prevent any congestion downline (the output line card applies retroactive rules to the input line cards to regulate their bit rate). This mechanism that should not hamper the flows with guaranteed service quality. Furthermore, since this mechanism tends to shift the downline congestion upline, it is up to the traffic management function, which is a sub-function of the ATM layer function, to reject cells that violate the policy upline storage resources allocation.

The French patent application No. 2 740 283 on behalf of the present Applicant presents an example of a flow control mechanism of this kind in Figure 1. With this basic mechanism, there are often associated complementary mechanisms for the elimination of "head-of-line blocking" based on specific queues in the input memory as a function of the outgoing junction. In practice, these mechanisms prove to be almost indispensable with regard to guaranteed service quality flows. Various examples of such mechanisms are presented and referred to in M. HYOJEONG SONG, "A simple and fast scheduler for input queued ATM switches", in HPC ASIA 97.

Symbolically, the point-to-multipoint connections may be represented by a "tree" with a "root" representing the sender user and its "leaves" representing the receiver users. The implementation of connections of this type is standardized as regards signalling and the routing. What has to be done simply is to form a point-to-multipoint connection by creating first of all a point-to-point connection and then grafting new leaves on it. This addition of leaves may be done at the initiative of the root or else of the leaf.

With regard to the conveyance of the cells, the point-to-point connection model, namely the input translation alone, cannot always be applied. Figure 3a, where the elements similar to those of Figure 1b are identified with the same references, shows a point-to-multipoint connection. This connection enters the switch by a port P1 and exits therefrom by the ports P2, P4, P5, P7. In this case, the input translation envisaged further above may order the switch fabric 5 to copy each cell of this connection into the three switch fabric junctions ( $7_3$ ,  $7_4$ ,  $7_n$ ) concerned but it is not capable of indicating the output ports to which the cell must be sent. To do so, it is

整理番号 P - 8 1 1 7

- 1 1 -

necessary to add this information into the translation tables and send it from the input up to the output. Furthermore, the output logic channel depends on each output port and it is not possible to assign a single logic channel to all the outputs.

For these reasons, it is generally preferred to make a twofold translation: an input translation which replaces the logic channel of the cell by a « broadcasting index » representing the connection within the switch, and then a second output translation which expresses the broadcasting index into a list of pairs of the form 'port, logic channel' with the copies that may be necessary.

With regard to the management of resources for the non-real-time flows, on the contrary, it is difficult to hope for a simple extension of mechanisms for the control of flows implemented for the point-to-point connections. Indeed, the head-of-line blocking mechanisms, making use of queues specific to each output, would dictate the management of as many queues as there are possible sets of outputs, namely  $2^N - 1$  in the case of N switch fabric junctions.

Another solution may consist, for example, in making as many copies of each cell in the input line card as there are output line cards concerned by this cell. This approach would have the unacceptable drawback of requiring N switch fabric cycles to transmit a cell to N output junctions, while the technology of the switch fabric generally enables the transfer of the cell towards N outputs in a single cycle.

By contrast, the French patent application No. 2 740 283 referred here above describes a single solution based on the local conversion of the different point-to-multipoint UBR/ABR flows into a single point-to-multipoint CBR flow by separating the transmissions and letting them through one by one from the input line card to the switch fabric and by means of an outright fixed-value reservation of the resources on all the output ports, the reservation level being fixed according to the bit rate dictated by a separating device known as a shaper.

However, this approach has drawbacks when the output ports have bit rates that are very different from one another. This is easily the case in an ATM switch where 25 Mbit/s ports can coexist with 622 Mbit/s ports. Indeed, in this case, it will be necessary to adjust the bit rate of the

shaper on the basis of the bit rate of the port with the smallest bit rate. Since all the ABR/UBR connections use the same shaper, the result thereof is a choking of all the connections, even of the connections that take only high bit rate output ports.

The multipoint-to-point and multipoint-to-multipoint connections are not presently being processed in the standardization units that deal with the ATM. There is therefore, for the time being, no signalling or routing defined for this type of connection.

In terms of the conveyance of the cells, any topology of communication leading data elements of different geographical origins to converge on one and the same link raises what is called the problem of the 'interlacing of PDU (protocol data unit) application frames'. Indeed, since the PDU application frames are segmented by the AAL layer into cells, the cells of different frames reach the addressee in an interlaced form. To reassemble the frames, the addressee should be capable of finding the frame to which each cell belongs. Now the segmenting mechanism most commonly used in the UBR connections, implemented in the AAL 5 adaptation layer, do not enable this identification. It enables only the identification of the last cell of the PDU frame, which is sufficient in point-to-point or point-to-multipoint modes for the ATM cells are transmitted in sequences.

For the switches with centralized architecture or with storage at output, it is possible to be satisfied with an approach that consists in keeping the cells of one and the same PDU frame in the memory, and in sending them all together when the PDU frame has been completely received. This is not the case for switches with input storage for the simultaneous access to the switch fabric shared by the different input line cards necessarily introduces another interlacing of the if the transmission from the different line cards is not coordinated.

Despite all these problems, there are existing needs of communications in multipoint-to-point and multipoint-to-multipoint modes. They could be dealt with in theory by the superimposition of point-to-point communications or point-to-multipoint communications. This case has appeared especially in the context of the emulation of local area networks known as LAN Emulation or LANE (local area network emulation) where any

sophisticated mechanism is installed to emulate a shared medium (ELAN) within an ATM network [AF LANE]. An example that may be referred to is that of a server known as a BUS (broadcast or unknown server) which is defined in the LANE standard. This server, whose architecture is shown in Figure 3b, enables a user to broadcast messages to all the users of an emulated shared medium or towards another user with whom he is not directly connected. Each user of the ELAN has a point-to-point connection towards the BUS server and the BUS server has a point-to-multipoint connection towards all the users of the ELAN, as indicated in Figure 3b. In the prior art, the broadcasting servers of the BUS server type are implemented either by an expensive dedicated piece of hardware connected somewhere to the ATM network or by a specific process that is executed on the management module of the switch. This last-named approach unfortunately is not greatly extensible when the number of users of the ELAN increases. Indeed, during the period of time when a new user is not known to his partners, the messages that he sends and receives must go through the BUS server. This is often a cause of congestion that may be fatal to the management module.

Another example of a need for multipoint-to-point communications and multipoint-to-multipoint connections is given by the emulation of routing between local area networks. This function in particular may be implemented according to the norm known as MPOA (multiprotocol over ATM) of the ATM Forum which makes it possible to obtain virtual routing between different emulated local networks (ELANs) or different virtual local area networks (VLANs) of an ELAN [AF MPOA]. Another way of obtaining the emulation of the routing consists in installing a routing software program in the management unit of the ATM switch. In this context, a virtual router can be likened a user of the different ELANs that it interconnects. On this basis, it must house a function, known as a LEC (LAN emulation client) for each ELAN. The virtual router must be implemented in a switch, for example by a specific process that is executed in the management module. This entails the risk of causing congestion in said module when the exchanges between the ELANs reach a sufficiently sustained level.

#### SUMMARY OF THE INVENTION

整理番号 P - 8 1 1 7

- 1 4 -

The aim of the invention is to overcome the above-mentioned drawbacks.

To this end, an object of the invention is a method for the control of flows within an ATM switch with distributed architecture and storage at input, wherein the method consists of the distribution, to each input line card, of a specified number  $n$  of shapers, a shaper being dedicated to the VBR<sub>nrt</sub> (variable bit rate non-real time) category flows and being adjusted as a function of the totalized mean bit rate and the  $n-1$  other shapers being dedicated to the ABR and UBR category flows and being adjusted as a function of the available bit rate (AvCR) of the output ports.

Other characteristics and advantages of the invention shall appear from the following description made with reference to the appended drawings.

#### MORE DETAILED DESCRIPTION

The method according to the invention overcomes the drawbacks of the prior art referred to here above in the case of point-to-multipoint connections. It does so by the implementation on each input line card 7, in particular as shown in Figure 4, no longer of one shaper but of several shapers  $9_1$  to  $9_n$ . More specifically, a shaper is dedicated to real-time flows of the category VBR<sub>nrt</sub>. The bit rate of this shaper is adjusted as a function of the totalized mean bit rate of the VBR<sub>nrt</sub> flows. A fixed number  $n-1$  of shapers dedicated to the UBR/ABR flow is adjusted as a function of the available rate (AvCR) of the output ports.

A range of operation is assigned to each of the UBR/ABR shapers identified by its number.

The ranges of operation characterize the bit rate values that a shaper may achieve. They are bit rate intervals (for example: 8 to 32 Mbit/s). The ranges of operation of the different shapers  $9_1$  to  $9_n$  are assigned definitively when the system is initialized and in such a way that the entire range of the nominal bit rates of the output ports is covered. The range of operation of a  $i+1$  ranking shaper has a lower limit lower than the lower limit of the range of the  $i$  ranking shaper and an upper limit slightly higher than the lower limit of the range of the  $i$  ranking shaper so as to obtain a slight overlapping of the range of operation of the adjacent shapers.

整理番号 P - 8 1 1 7

- 1 5 -

For example, it is possible to envisage the following distribution in a switch whose output ports have nominal bit rates which are stepped in intervals from 64 kbit/s to 155 Mbit/s:

number	range of operation
shaper 1	155 Mbit/s - 8 Mbit/s
shaper 2	32 Mbit/s - 2 Mbit/s
shaper 3	8 Mbit/s - 64 kbit/s

With each point-to-multipoint connection, there is associated a "tree" constituted by the data element of its input line card (root) and its output ports (leaves). (This terminology should not be confused with the previous terminology used in the prior art description where "root" and "leaf" are applied to users.) As and when the life of the connection progresses, the tree that is associated with it may get modified by the addition or elimination of leaves.

According to the invention, with each tree, there is associated a shaper which processes the cells of all the point-to-multipoint connections associated with this tree. The choice of the shaper is made as a function of the nominal bit rate of the "leaf" port having the lowest nominal bit rate, with a law of assignment that is within the scope of those skilled in the art. In the context of the above example, it may be decided that:

- the trees for which all the "leaf" ports have a nominal bit rate of 155 Mbit/s shall be assigned to the shaper 1,
- the trees that do not come within the above category and of which no "leaf" port has a nominal bit rate, below 25 Mbit/s, will be assigned to the shaper 2,
- the trees that do not come under the above two categories will be assigned to the shaper 3.

In this way, each point-to-multipoint connection is assigned to a shaper, this assignment being revisable at each modification of the tree of the connection.

While the assignment of the connections to the shapers  $9_1$  to  $9_n$  is done as a function of the nominal bit rate of the leaves, the adjusting of the bit rate of the switches, on the contrary, is done as a function of the available bit rate of the leaves. The available bit rate is a difference between the



nominal bit rate and the reserved bit rate. The concept of the reserved bit rate is one that depends on the implementation, which represents the totalizing of the bit rates guaranteed for the connections CBR, VBR, ABR and even UBR if the policy of the network means ensuring a minimum bit rate for the UBR connections. The bit rate of a shaper is computed by taking the minimum of the available bit rates of the leaves belonging to trees assigned to the shaper concerned. The precise implementation of a mechanism of this kind for adjusting the bit rates of shapers is within the scope of those skilled in the art. In the context of the above example, it may be assumed for example that the management module of the switch keeps the available bit rate of each output port up to date and that it periodically sends each input line card the bit rate to be applied to each of its shapers.

With regard to the shapers dedicated to the VBRrt real time flows, these are adjusted as a function of the totalized mean bit rate of the connections concerned.

By the application of the above mechanism, following the major drop in the available bit rate on one or more output ports, the bit rate allocated to a rank  $i$  shaper could be incompatible with its range of operation. Thus, to avoid this situation, it is planned to have an exceptional mechanism by means of which those trees whose leaves are the ports concerned will be temporarily reassigned to a shaper adapted to lower bit rates, for example the  $i+1$  rank shaper. Here the overlapping of the ranges of operation creates a natural hysteresis effect as a result of which the changes of assignment do not take place with excessive frequency and prevent excessively frequent changes of shapers for one and the same tree.

All the ABR/UBR shapers having the same number, located on different input line cards, are capable of sending cells simultaneously to one and the same output port, each at a bit rate lower than or equal to the available bit rate of this port. The result of this is that, in the worst case, the bit rate coming from the input line cards can easily go beyond the available bit rate of the concerned port. This problem is resolved by the implementation of a mechanism for the arbitration of access by the different line cards to the scrambling resource. This mechanism also makes it possible to overcome the prior art drawback referred to, concerning the interleaving of the PDU application frames due to the AAL5 segmentation in

the context of multipoint-to-point or multipoint-to-multipoint connections.

To do this, each of the ABR/UBR connection shapers arbitrates between two queues: a priority queue P1 which works in PDU application frame mode and a queue P2 which works in cell mode. To say that a queue works in PDU mode means that it is not ready to transmit unless it contains at least one full frame and unless the cells of one and the same frame are transmitted sequentially without any possibility for cells of other frames to get interposed. To say that a queue works in cell mode means that it is ready to transmit as soon as it contains at least one cell.

Then, on each input line card, the multipoint-to-point connections and multipoint-to-multipoint connections are assigned to trees, in a manner similar to that described for the point-to-multipoint connections. For the multipoint-to-point connections, the tree is reduced to its root and to a single leaf. For the multipoint-to-multipoint connections, the tree may take, as its root, the input line card considered and, as its leaves, all the ports taking part in the connection. The trees will be assigned to the shapers according to the same principles as those stipulated here above.

The multipoint-to-point connections and multipoint-to-multipoint connections use only the queues P1 for they run the risk of interlacing of the AAL5 PDUs.

Each of the ABR/UBR shapers  $9_1$  to  $9_n$  of each line card  $7_i$  is assigned to a class. There is a distribution into classes, namely a division of all the shapers of a given rank  $i$ . This division is different for each rank of shaper, namely for each range of operation. The classes are defined by the following principle: in a given range of operation, a shaper of this range is in the same class as another shaper of this range if at least one tree associated with the first shaper and at least one tree associated with the second shaper has at least one common leaf. Conversely, two shapers are in different classes, as can be seen in Figure 5a, if all the leaves of the trees associated with one of them has a vacant intersection with all the leaves of the trees associated with the other shaper. It must be noted that this division is dynamic: the assignment of a new tree to a shaper or its elimination will generally require the recomputation of the division. Each division comprises at least one class and, at most, as many classes as there are line cards in the switch.

The algorithm for the classification of the shapers according to the principle stipulated here above is within the scope of those skilled in the art.

The principle of arbitration is based up a division of time into intervals with a length  $T = T1 + T2 + T3$ . A « order of speaking » is assigned to the various shapers of the same class. At each time interval  $T$ , the shaper of each class that has the « right to speak » can begin sending the cells of a PDU extracted from the file  $P1$  during the interval  $T1$ . During the interval  $T2$ , it may continue to send the cells of a PDU whose transmission has already begun but it cannot start sending a new PDU. This same shaper may send cells coming from  $P2$  throughout the duration of the interval  $T1 + T2$ , provided of course that it is incapable of sending cells coming from  $P1$  (according to the principle of priority of  $P1$  over  $P2$ ). If its transmission ends before the end of  $T1 + T2$ , a shaper may hand over the right to speak to another shaper of its class which itself could make transmission during the same period. Once the interval  $T1 + T2$  has elapsed, each line card indicates its desire to send a message by means of a collecting mechanism that runs in the interval  $T3$ . The « turn to speak » during the interval  $T$  will also be computed as a function of wishes expressed and in favoring a shaper that has had to suspend the transmission of a PDU before the end of this PDU. The algorithm for computing the « turn to speak » must ensure a certain fairness between the line cards if the management policy of the switch requires it. An algorithm of this kind is within the scope of those skilled in the art. All that is said here above pertains to a given rank  $i$  of a shaper. Thus, for example, the intervals  $T$ ,  $T1$ ,  $T2$ ,  $T3$  depend on  $i$ , namely on the range of operations of the shaper considered.

An implementation of the above principles based on tokens is described hereinafter by way of an example. Different phases of this mechanism shown in Figures 5b, 5c and 5d are described.

The tokens have specialized cells exchanged with the different line cards  $7_i$ , on the initiative of a master unit which may be one of the line cards  $7_i$  or else the management module 4 or else the switch fabric 2 itself. The invention distinguishes three types of token.

The « SYNC SIGNAL token » of Figure 5b is sent periodically by the master unit with a periodicity  $T$ , in a general broadcast to all the line cards whose  $i$  rank shaper possesses at least one tree that is associated

with it. This token is used for the approximate synchronizing of the transmissions between the different line cards. It also conveys the « turn to speak » for each class.

The « SPEECH token » of Figure 5c is specific to each of the classes identified here above. It is sent by one of the line cards whose rank i shaper has finished the transmission to the next line card, before the end of the interval  $T1 + T2$ , of the « turn to speak » within the class considered.

The « COLLECTION token » of Figure 5d is sent by one of the line cards that has received the SYNC SIGNAL token at the end of the period T. It is relayed from one unit to the next in all the concerned line cards which record an indicator therein on their need to make transmission during the following period or to continue to transmit the end of a PDU whose transmission has already started. Finally, the token is received and processed by the master unit.

To overcome the prior art drawback concerning the case of broadcasting servers, such as BUS servers, made by the superimposition of point-to-point connections of each user of a broadcasting group towards the server and a point-to-multipoint connection of the server towards all the users of the same group, a short-circuit is made within the ATM switch between the different connections. This enables the creation of a multipoint-to-multipoint type communication in which the relaying of the PDU application frames is done totally by the junction devices, the server no longer having any role other than that of the management of the broadcasting groups, namely the operations needed for the arrival of a new user in the group or the departure of a user from a group. The process of relaying PDU frames at the level of the line cards is made possible by the implementation of the arbitration and flow control processes described here above.

The broadcasting server which bears the reference 10 in Figure 6a is a software process that is executed on the management module 4 of the switch. In the prior art mode of operation, the communications module is the one shown in Figure 6b which illustrates the example of three users belonging to one broadcasting group. In the mode of operation according to the invention, the communications module is the one shown in Figure 6c in which the point-to-point connections A1, A2, A3 and the point-to-multipoint

整理番号 P - 8 1 1 7

- 2 0 -

connection M are replaced by point-to-multipoint connections B1, B2, B3 from each user of the group towards all the others.

This modification of the communications module is made possible by a simple conversion of the input and output translations implemented in the ATM layer and referenced TRANS IN, TRANS OUT in the line cards 7<sub>1</sub>, 7<sub>2</sub>, 7<sub>3</sub> of Figure 6a. With regard to the input translations of the Ai type point-to-point connections, they are, in the prior art.  $VLAi \rightarrow (VLA'i, LG)$  where VLAi is a logic path identifier associated with Ai on the input artery, VLA'i is the identifier associated with this same connection in the management module and LG is the identity of the switch fabric junction of the management module. This expresses the fact that all the PDU frames are processed by the server in the implementation module. Conversely, all the frames retransmitted by the server to the users of a given broadcasting group are sent with a logic channel VLM characteristic of the point-to-multipoint connection M. An output translation makes it possible, in each output line card concerned, to convert the broadcasting index VLM into a logic channel identifier VLMI characteristic of this connection on the output artery with a possible copying of each cell if it must exit at several ports P1, P2, ... Pp. This last-named translation therefore has the form  $VLM \rightarrow (VLMP1, P1), (VLMP2, P2), \dots, (VLMPp, Pp)$ .

According to the invention, the only translations modified by the short-circuit are the input translations which become  $VLAi \rightarrow (VLM, L1, L2, L3, \dots, Lq)$  where L1 to Lq designate all the switch fabric junctions concerned by the connection M.

In the mode of operation according to the invention, the broadcasting server manages the arrival and departure of the users in the same way as in the prior art and then determines the input translations to be modified in order to arrive at the communications module described here above.

The table of Figure 6d shows an exemplary translation with and without short-circuits in the case of a broadcasting group comprising four users ui, u21, u22 and u3 as indicated in Figure 6a.

Note 1: without being exhaustive, Figure 6b represents a case where the connections A1, A2, A3 are one-way connections which is not the case for the LANE server BUS.

整理番号 P - 8 1 1 7

- 2 1 -

Note 2: the connections A1, A2, A3 and M continue to exist at least with regard to the signalling.

Note 3: depending on the applications, it is always possible that the server itself may or may not be considered to form part of the group of users.

Finally, in addition to the arbitrating mechanism described here above, it is possible to have a third mechanism that makes it possible to overcome the prior art drawback referred to here above if a frame router IP is implemented in a management module. The method makes it possible, within the ATM switch, to obtain a real decentralization of the relay function IP by limiting the role of the router to the computation of the routes (prior art).

Figure 3c shows a case where this method can be used. The internal router has as many LECs as there are ELANs known to it. If the user uA belonging to the ELAN A wishes to send a frame IP to the user uB, he starts by using the broadcasting means on the emulator ELAN (BUS broadcasting server). If the router function internal to the switch knows about the existence of the user uB, the function LEC A of this router, associated with the ELAN A emulation, declares itself to be the addressee of all the IP frames addressed to the user uB. Thereafter, the user uA uses his direct ATM connection to perform the LEC A function, according to usual procedures specified in the LANE standard for sending frames to the user uB. In the prior art, the frames must go back up to the internal router at the management module in order to be relayed towards the user uB through the use of the direct ATM connection existing between the LEC B function and the user uB.

The method according to the invention specifies that, for any PDU application frame arriving at a direct connection relating to the LEC A function, the input line card examines its first cell and extracts the IP address of the addressee therefrom. It then travels through a cache table updated by means of routing information coming from the management module and there, before the IP address, it finds a pair (logic channel, outgoing direction). The logic channel is the one used by the direct connection between the LEC B and uB. The outgoing direction is the identifier of the switch fabric junction concerned by this direction connection. If the input line card does not find the IP address sought in the cache table, it sends the

management module a request for updating the cache. The information found in the table is then used for the translation of the ATM header of each cell of the PDU concerned. Thus, a dynamic translation is obtained, the translation table being modified potentially during the passage of each PDU. Thus, a "dynamic short-circuit" is set up between two point-to-point connections.

#### 4 Brief Description of Drawings

- Figure 1a is a drawing showing the principle of an ATM switch according to the invention;
- Figure 1b is a drawing showing the principle of an ATM switch with distributed architecture according to the prior art;
- Figures 2a to 2f are drawings showing modes of communication between users of an ATM network;
- Figure 3a shows an example of a routing of an ATM cell in a switch during a point-to-multipoint connection;
- Figure 3b shows an example of the superimposing of point-to-point connections or point-to-multipoint connections in an emulated LAN architecture;
- Figure 3c is a drawing showing the principle of routing between ELANs;
- Figure 3d shows a principle of short-circuit implemented by the invention;
- Figure 4 is a drawing showing the principle of the assignment of trees to shapers according to the invention;
- Figure 5a shows an example illustrating an operation of the classifying of shapers according to the invention;
- Figures 5b, 5c, and 5d show examples of the broadcasting of tokens according to the invention;
- Figures 6a to 6d show examples of implementation of a broadcasting server according to the invention.

FIG. 1 a

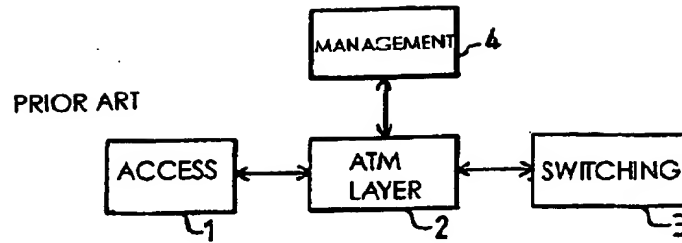


FIG. 1 b

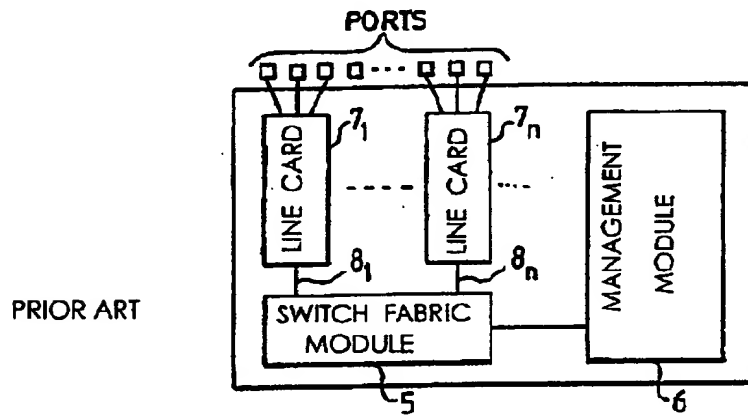


FIG. 2 a

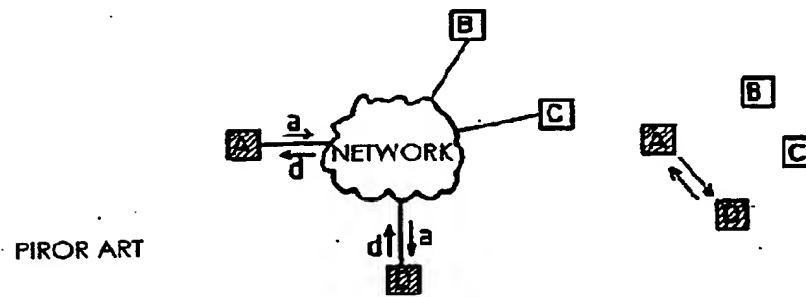
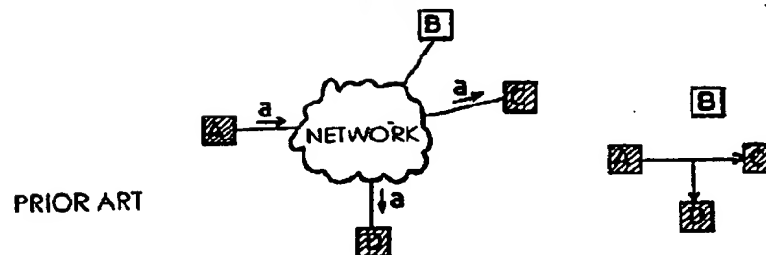


FIG. 2 b





整理番号 P-8117

- 2 -

FIG. 2c

PRIOR ART

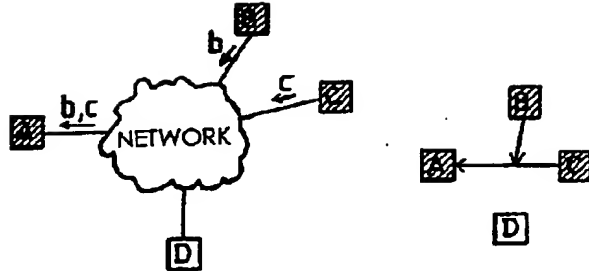


FIG. 2d

PRIOR ART

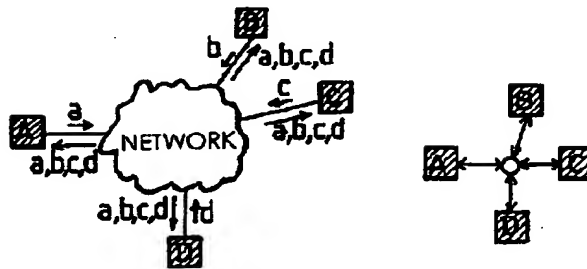


FIG. 2e

PRIOR ART

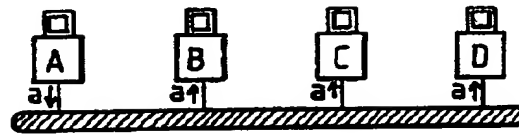
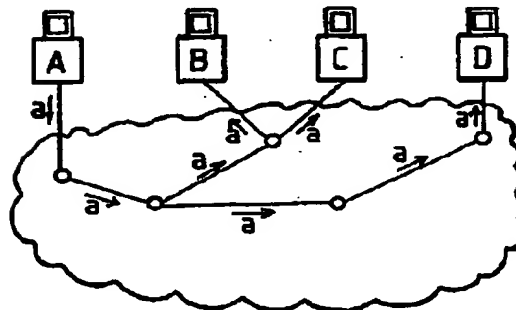


FIG. 2f

PRIOR ART



整理番号 P-8117

- 3 -

FIG. 3a

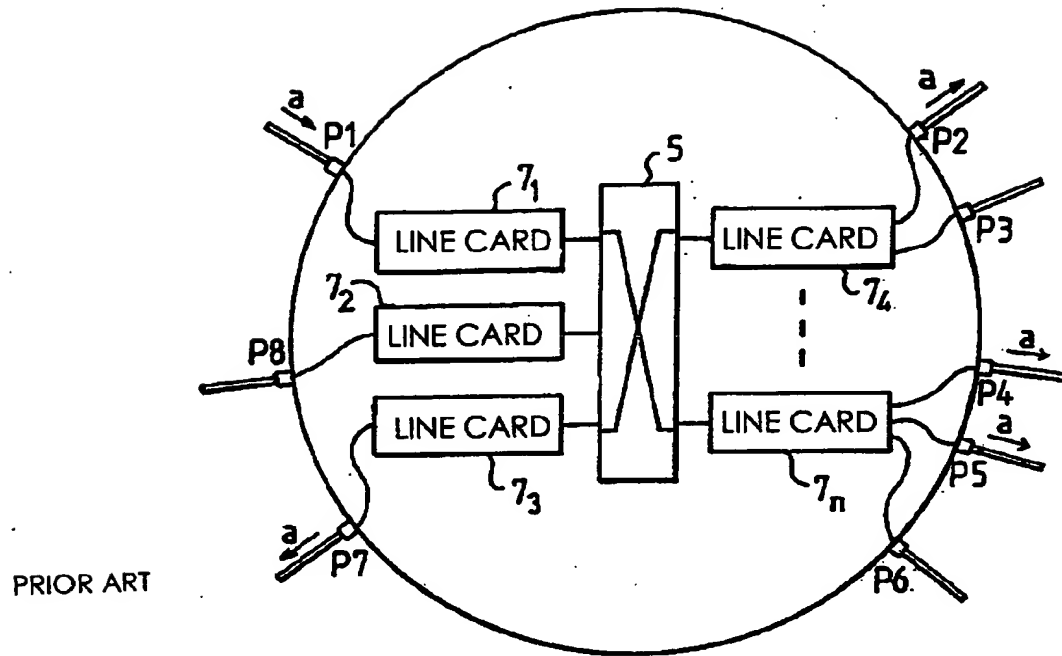
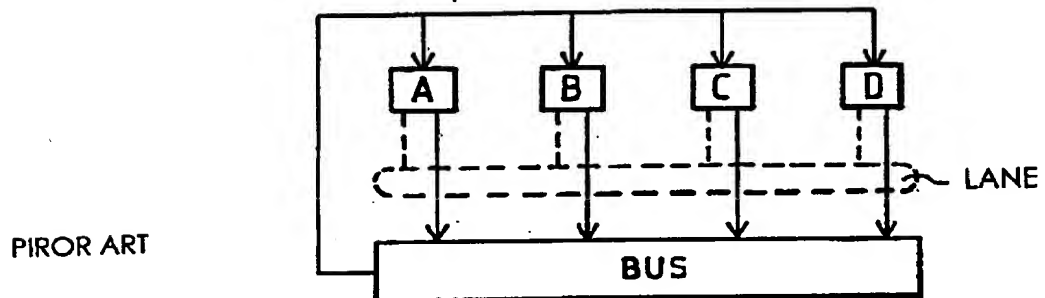


FIG. 3b



整理番号 P-8117

- 4 -

FIG. 3c

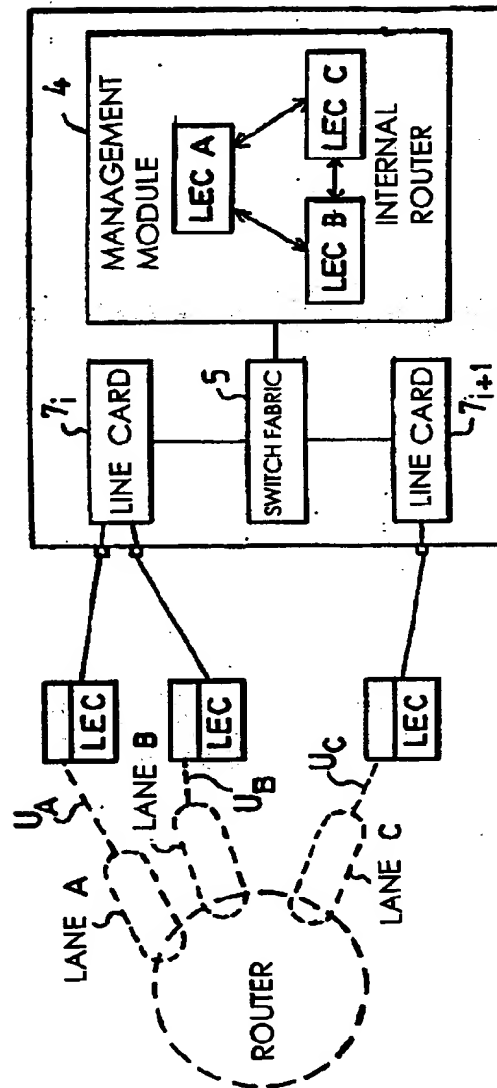


FIG. 3d

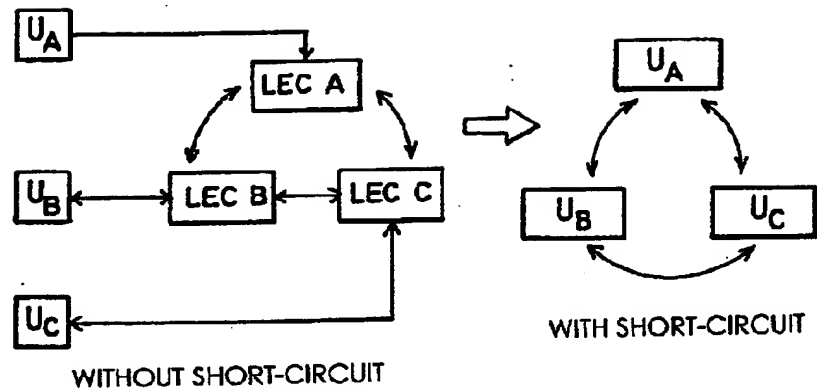


FIG. 4

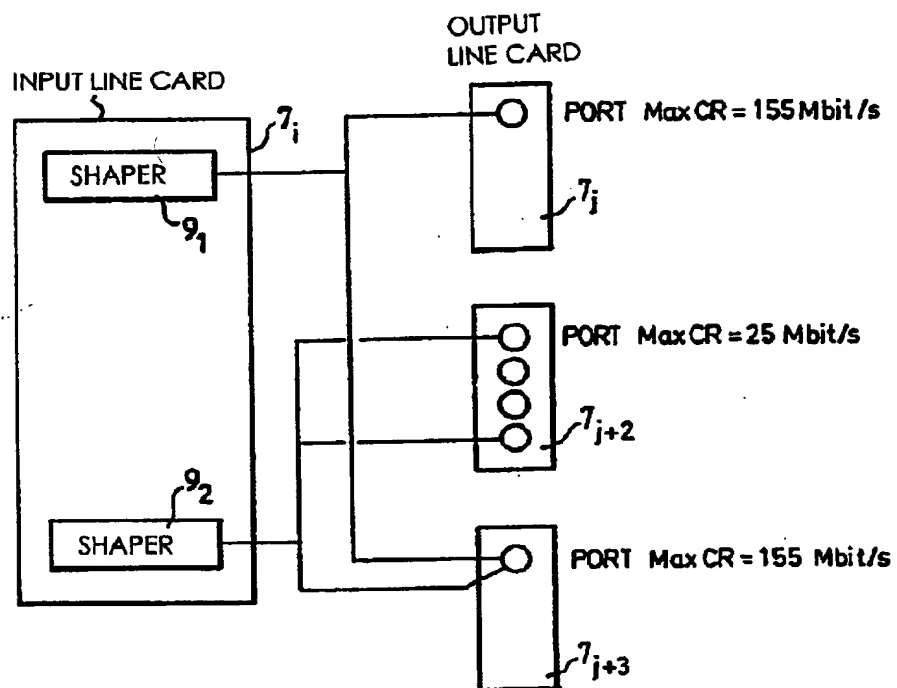


FIG. 5 a

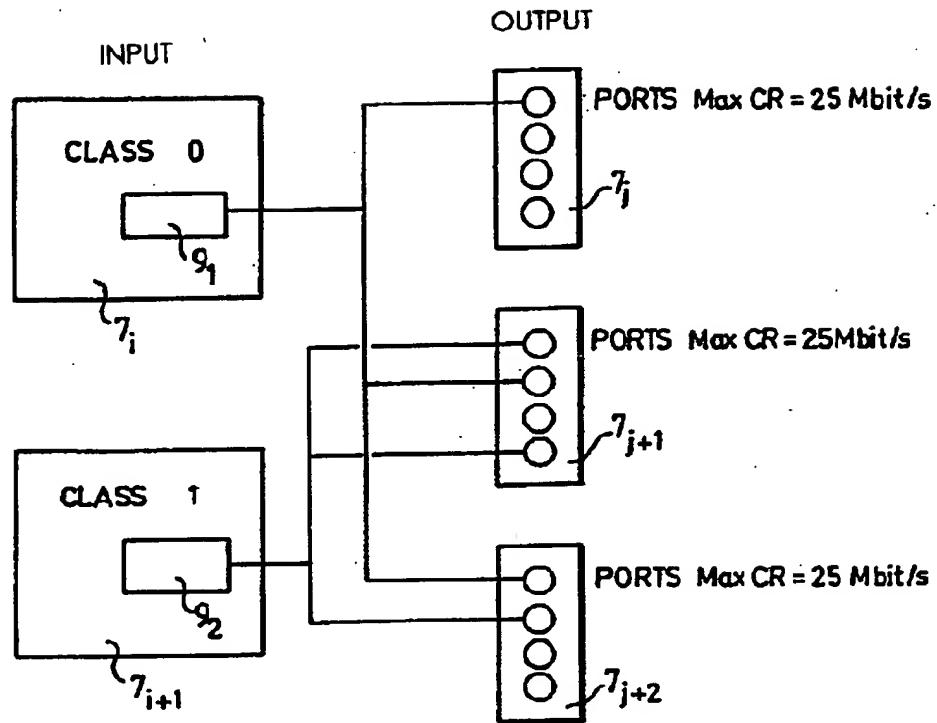
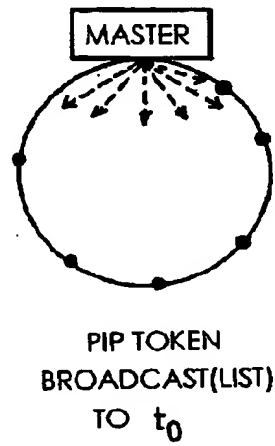


FIG. 5 b



整理番号 P-8117

- 7 -

FIG. 5c

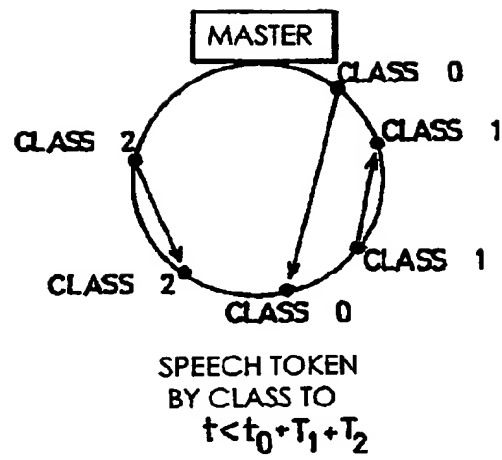
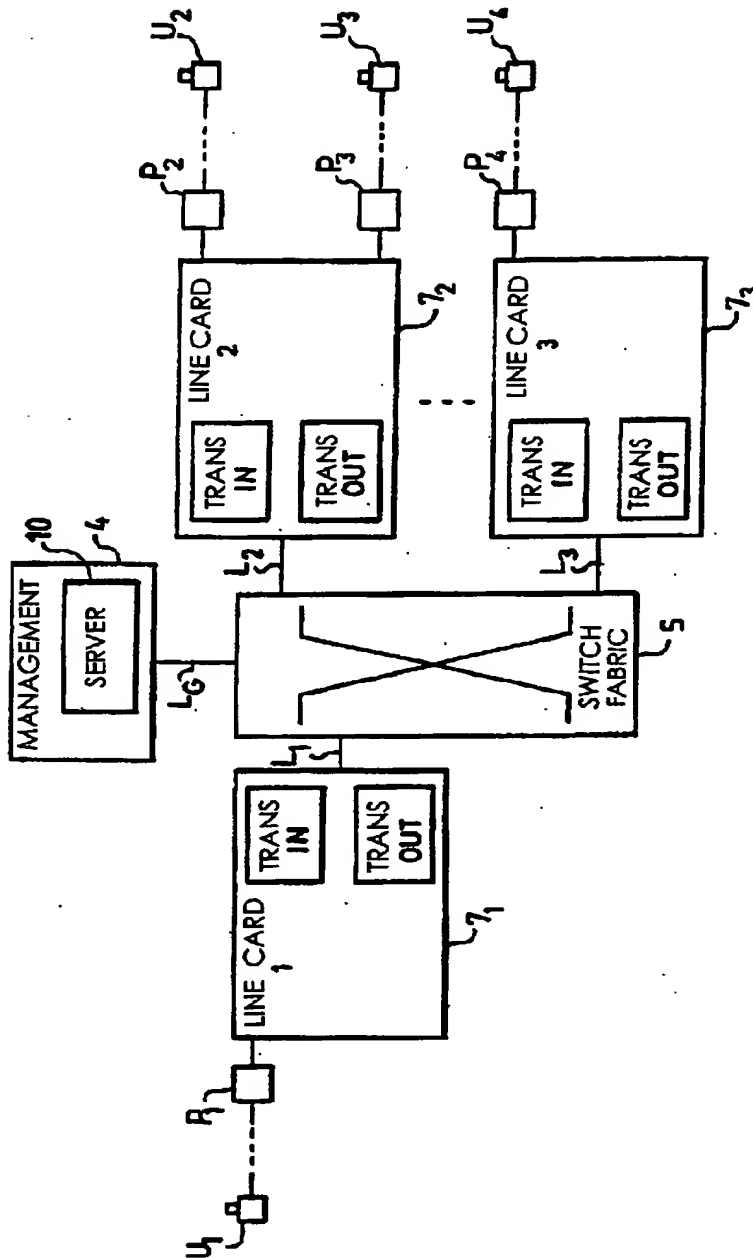


FIG. 5d



FIG. 6a



整理番号 P-8117

- 9 -

FIG. 6b

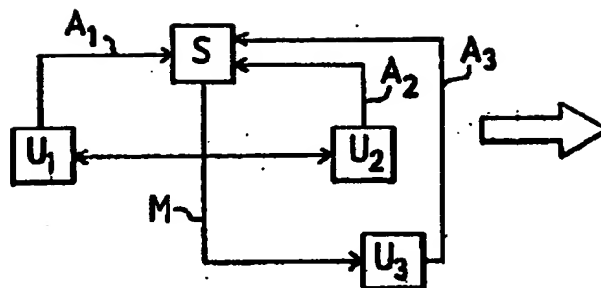
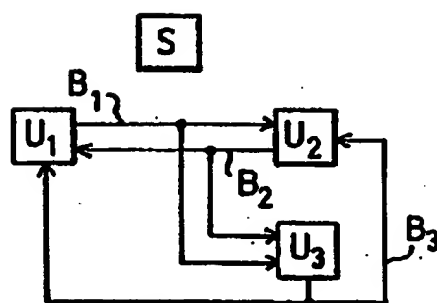


FIG. 6c





整理番号 P-8117

- 10 -

FIG. 6d

	WITHOUT SHORT-CIRCUIT	WITH SHORT-CIRCUIT
LINE CARD 1, INPUT	VLA 1 → VLA'1, LG	VLA 1 → VLM, L1, L2, L3
LINE CARD 1, OUTPUT	VLM → VLM1, P1	DITTO
LINE CARD 2, INPUT	VLA 21 → VLA'21, LG VLA 22 → VLA'22, LG	VLA 21 → VLM, L1, L2, L3 VLA 22 → VLM, L1, L2, L3
LINE CARD 2, OUTPUT	VLM → (VLM 21, P21) (VLM 22, P22)	DITTO
LINE CARD 3, INPUT	VLA 3 → VLA'3, LG	VLA 3 → VLM, L1, L2, L3
LINE CARD 3, OUTPUT	VLM → VLM3, P3	DITTO

## 1 Abstract

The method consists in distributing a specified number  $n$  of shapers ( $g_1 \dots g_n$ ) to each input line card ( $7_i$ ), a shaper being dedicated to the VBRnrt (Variable Bit Rate non real time) category flows as a function of the totalized mean bit rate and the  $n-1$  other shapers being adjusted as a function of the available bit rate (AVCR) of the output ports. Application to ATM transmission networks..

## 2 Representative Drawing

FIGURE 4